

# Algorithms and Computation in Mathematics • Volume 10

*Editors*

Arjeh M. Cohen    Henri Cohen  
David Eisenbud    Michael F. Singer  
Bernd Sturmfels

Saugata Basu  
Richard Pollack  
Marie-Françoise Roy

# Algorithms in Real Algebraic Geometry

Second Edition

With 37 Figures

 Springer

Saugata Basu  
Georgia Institute of Technology  
School of Mathematics  
Atlanta, GA 30332-0160  
USA  
e-mail: saugata@math.gatech.edu

Richard Pollack  
Courant Institute of  
Mathematical Sciences  
251 Mercer Street  
New York, NY 10012  
USA  
e-mail: pollack@cims.nyu.edu

Marie-Françoise Roy  
IRMAR Campus de Beaulieu  
Université de Rennes I  
35042 Rennes cedex  
France  
e-mail: Marie-Francoise.Roy@univ-rennes1.fr

Library of Congress Control Number: 2006927110

---

Mathematics Subject Classification (2000): 14P10, 68W30, 03C10, 68Q25, 52C45

---

ISSN 1431-1550

ISBN-10 3-540-33098-4 Springer Berlin Heidelberg New York

ISBN-13 978-3-540-33098-1 Springer Berlin Heidelberg New York

ISBN 3-540-00973-6 1st edition Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable for prosecution under the German Copyright Law.

Springer is a part of Springer Science+Business Media  
springer.com

© Springer-Verlag Berlin Heidelberg 2003, 2006  
Printed in Germany

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typeset by the authors using a Springer  $\text{\LaTeX}$  macro package  
Production: LE- $\text{\TeX}$  Jelonek, Schmidt & Vöckler GbR, Leipzig  
Cover design: *design & production* GmbH, Heidelberg

Printed on acid-free paper 46/3100YL - 5 4 3 2 1 0

---

# Table of Contents

<b>Introduction</b> . . . . .	1
<b>1 Algebraically Closed Fields</b> . . . . .	11
1.1 Definitions and First Properties . . . . .	11
1.2 Euclidean Division and Greatest Common Divisor . . . . .	14
1.3 Projection Theorem for Constructible Sets . . . . .	20
1.4 Quantifier Elimination and the Transfer Principle . . . . .	25
1.5 Bibliographical Notes . . . . .	27
<b>2 Real Closed Fields</b> . . . . .	29
2.1 Ordered, Real and Real Closed Fields . . . . .	29
2.2 Real Root Counting . . . . .	44
2.2.1 Descartes's Law of Signs and the Budan-Fourier Theorem . . . . .	44
2.2.2 Sturm's Theorem and the Cauchy Index . . . . .	52
2.3 Projection Theorem for Algebraic Sets . . . . .	57
2.4 Projection Theorem for Semi-Algebraic Sets . . . . .	63
2.5 Applications . . . . .	69
2.5.1 Quantifier Elimination and the Transfer Principle . . . . .	69
2.5.2 Semi-Algebraic Functions . . . . .	71
2.5.3 Extension of Semi-Algebraic Sets and Functions . . . . .	72
2.6 Puiseux Series . . . . .	74
2.7 Bibliographical Notes . . . . .	81
<b>3 Semi-Algebraic Sets</b> . . . . .	83
3.1 Topology . . . . .	83
3.2 Semi-algebraically Connected Sets . . . . .	86
3.3 Semi-algebraic Germs . . . . .	87

3.4	Closed and Bounded Semi-algebraic Sets	93
3.5	Implicit Function Theorem	94
3.6	Bibliographical Notes	99
<b>4</b>	<b>Algebra</b>	101
4.1	Discriminant and Subdiscriminant	101
4.2	Resultant and Subresultant Coefficients	105
4.2.1	Resultant	105
4.2.2	Subresultant Coefficients	110
4.2.3	Subresultant Coefficients and Cauchy Index	113
4.3	Quadratic Forms and Root Counting	119
4.3.1	Quadratic Forms	119
4.3.2	Hermite's Quadratic Form	127
4.4	Polynomial Ideals	132
4.4.1	Hilbert's Basis Theorem	132
4.4.2	Hilbert's Nullstellensatz	136
4.5	Zero-dimensional Systems	143
4.6	Multivariate Hermite's Quadratic Form	149
4.7	Projective Space and a Weak Bézout's Theorem	153
4.8	Bibliographical Notes	157
<b>5</b>	<b>Decomposition of Semi-Algebraic Sets</b>	159
5.1	Cylindrical Decomposition	159
5.2	Semi-algebraically Connected Components	168
5.3	Dimension	170
5.4	Semi-algebraic Description of Cells	172
5.5	Stratification	174
5.6	Simplicial Complexes	181
5.7	Triangulation	183
5.8	Hardt's Triviality Theorem and Consequences	186
5.9	Semi-algebraic Sard's Theorem	191
5.10	Bibliographical Notes	194
<b>6</b>	<b>Elements of Topology</b>	195
6.1	Simplicial Homology Theory	195
6.1.1	The Homology Groups of a Simplicial Complex	195
6.1.2	Simplicial Cohomology Theory	199
6.1.3	A Characterization of $H^1$ in a Special Case.	201
6.1.4	The Mayer-Vietoris Theorem	206

6.1.5	Chain Homotopy	209
6.1.6	The Simplicial Homology Groups Are Invariant Under Homeomorphism	213
6.2	Simplicial Homology of Closed and Bounded Semi-algebraic Sets	221
6.2.1	Definitions and First Properties	221
6.2.2	Homotopy	223
6.3	Homology of Certain Locally Closed Semi-Algebraic Sets	226
6.3.1	Homology of Closed Semi-algebraic Sets and of Sign Conditions	226
6.3.2	Homology of a Pair	228
6.3.3	Borel-Moore Homology	231
6.3.4	Euler-Poincaré Characteristic	234
6.4	Bibliographical Notes	236
<b>7</b>	<b>Quantitative Semi-algebraic Geometry</b>	<b>237</b>
7.1	Morse Theory	237
7.2	Sum of the Betti Numbers of Real Algebraic Sets	256
7.3	Bounding the Betti Numbers of Realizations of Sign Conditions	262
7.4	Sum of the Betti Numbers of Closed Semi-algebraic Sets	268
7.5	Sum of the Betti Numbers of Semi-algebraic Sets	273
7.6	Bibliographical Notes	280
<b>8</b>	<b>Complexity of Basic Algorithms</b>	<b>281</b>
8.1	Definition of Complexity	281
8.2	Linear Algebra	292
8.2.1	Size of Determinants	292
8.2.2	Evaluation of Determinants	294
8.2.3	Characteristic Polynomial	299
8.2.4	Signature of Quadratic Forms	300
8.3	Remainder Sequences and Subresultants	301
8.3.1	Remainder Sequences	301
8.3.2	Signed Subresultant Polynomials	303
8.3.3	Structure Theorem for Signed Subresultants	307
8.3.4	Size of Remainders and Subresultants	314
8.3.5	Specialization Properties of Subresultants	316
8.3.6	Subresultant Computation	317
8.4	Bibliographical Notes	322

<b>9</b>	<b>Cauchy Index and Applications</b>	323
9.1	Cauchy Index	323
9.1.1	Computing the Cauchy Index	323
9.1.2	Bezoutian and Cauchy Index	326
9.1.3	Signed Subresultant Sequence and Cauchy Index on an Interval	330
9.2	Hankel Matrices	333
9.2.1	Hankel Matrices and Rational Functions	334
9.2.2	Signature of Hankel Quadratic Forms	337
9.3	Number of Complex Roots with Negative Real Part	344
9.4	Bibliographical Notes	350
<b>10</b>	<b>Real Roots</b>	351
10.1	Bounds on Roots	351
10.2	Isolating Real Roots	360
10.3	Sign Determination	383
10.4	Roots in a Real Closed Field	397
10.5	Bibliographical Notes	401
<b>11</b>	<b>Cylindrical Decomposition Algorithm</b>	403
11.1	Computing the Cylindrical Decomposition	404
11.1.1	Outline of the Method	404
11.1.2	Details of the Lifting Phase	408
11.2	Decision Problem	415
11.3	Quantifier Elimination	423
11.4	Lower Bound for Quantifier Elimination	426
11.5	Computation of Stratifying Families	428
11.6	Topology of Curves	430
11.7	Restricted Elimination	440
11.8	Bibliographical Notes	444
<b>12</b>	<b>Polynomial System Solving</b>	445
12.1	A Few Results on Gröbner Bases	445
12.2	Multiplication Tables	451
12.3	Special Multiplication Table	456
12.4	Univariate Representation	462
12.5	Limits of the Solutions of a Polynomial System	471
12.6	Finding Points in Connected Components of Algebraic Sets	483
12.7	Triangular Sign Determination	495

12.8	Computing the Euler-Poincaré Characteristic of an Algebraic Set . . . . .	498
12.9	Bibliographical Notes . . . . .	503
<b>13</b>	<b>Existential Theory of the Reals . . . . .</b>	<b>505</b>
13.1	Finding Realizable Sign Conditions . . . . .	506
13.2	A Few Applications . . . . .	516
13.3	Sample Points on an Algebraic Set . . . . .	519
13.4	Computing the Euler-Poincaré Characteristic of Sign Conditions . . . . .	528
13.5	Bibliographical Notes . . . . .	532
<b>14</b>	<b>Quantifier Elimination . . . . .</b>	<b>533</b>
14.1	Algorithm for the General Decision Problem . . . . .	534
14.2	Quantifier Elimination . . . . .	547
14.3	Local Quantifier Elimination . . . . .	551
14.4	Global Optimization . . . . .	557
14.5	Dimension of Semi-algebraic Sets . . . . .	558
14.6	Bibliographical Notes . . . . .	562
<b>15</b>	<b>Computing Roadmaps and Connected Components of Algebraic Sets . . . . .</b>	<b>563</b>
15.1	Pseudo-critical Values and Connectedness . . . . .	564
15.2	Roadmap of an Algebraic Set . . . . .	568
15.3	Computing Connected Components of Algebraic Sets . . . . .	580
15.4	Bibliographical Notes . . . . .	592
<b>16</b>	<b>Computing Roadmaps and Connected Components of Semi-algebraic Sets . . . . .</b>	<b>593</b>
16.1	Special Values . . . . .	593
16.2	Uniform Roadmaps . . . . .	601
16.3	Computing Connected Components of Sign Conditions . . . . .	608
16.4	Computing Connected Components of a Semi-algebraic Set . . . . .	614
16.5	Roadmap Algorithm . . . . .	617
16.6	Computing the First Betti Number of Semi-algebraic Sets . . . . .	627
16.7	Bibliographical Notes . . . . .	633
	<b>References . . . . .</b>	<b>635</b>
	<b>Index of Notation . . . . .</b>	<b>645</b>
	<b>Index . . . . .</b>	<b>655</b>



---

## Introduction

Since a real univariate polynomial does not always have real roots, a very natural algorithmic problem, is to design a method to count the number of real roots of a given polynomial (and thus decide whether it has any). The “real root counting problem” plays a key role in nearly all the “algorithms in real algebraic geometry” studied in this book.

Much of mathematics is algorithmic, since the proofs of many theorems provide a finite procedure to answer some question or to calculate something. A classic example of this is the proof that any pair of real univariate polynomials  $(P, Q)$  have a greatest common divisor by giving a finite procedure for constructing the greatest common divisor of  $(P, Q)$ , namely the euclidean remainder sequence. However, different procedures to solve a given problem differ in how much calculation is required by each to solve that problem. To understand what is meant by “how much calculation is required”, one needs a fuller understanding of what an algorithm is and what is meant by its “complexity”. This will be discussed at the beginning of the second part of the book, in Chapter 8.

The first part of the book (Chapters 1 through 7) consists primarily of the mathematical background needed for the second part. Much of this background is already known and has appeared in various texts. Since these results come from many areas of mathematics such as geometry, algebra, topology and logic we thought it convenient to provide a self-contained, coherent exposition of these topics.

In Chapter 1 and Chapter 2, we study algebraically closed fields (such as the field of complex numbers  $\mathbb{C}$ ) and real closed fields (such as the field of real numbers  $\mathbb{R}$ ). The concept of a real closed field was first introduced by Artin and Schreier in the 1920’s and was used for their solution to Hilbert’s 17th problem [6, 7]. The consideration of abstract real closed fields rather than the field of real numbers in the study of algorithms in real algebraic geometry is not only intellectually challenging, it also plays an important role in several complexity results given in the second part of the book.

Chapters 1 and 2 describe an interplay between geometry and logic for algebraically closed fields and real closed fields. In Chapter 1, the basic geometric objects are constructible sets. These are the subsets of  $\mathbb{C}^n$  which are defined by a finite number of polynomial equations ( $P = 0$ ) and inequations ( $P \neq 0$ ). We prove that the projection of a constructible set is constructible. The proof is very elementary and uses nothing but a parametric version of the euclidean remainder sequence. In Chapter 2, the basic geometric objects are the semi-algebraic sets which constitute our main objects of interest in this book. These are the subsets of  $\mathbb{R}^n$  that are defined by a finite number of polynomial equations ( $P = 0$ ) and inequalities ( $P > 0$ ). We prove that the projection of a semi-algebraic set is semi-algebraic. The proof, though more complicated than that for the algebraically closed case, is still quite elementary. It is based on a parametric version of real root counting techniques developed in the nineteenth century by Sturm, which uses a clever modification of euclidean remainder sequence. The geometric statement “the projection of a semi-algebraic set is semi-algebraic” yields, after introducing the necessary terminology, the theorem of Tarski that “the theory of real closed fields admits quantifier elimination.” A consequence of this last result is the decidability of elementary algebra and geometry, which was Tarski’s initial motivation. In particular whether there exist real solutions to a finite set of polynomial equations and inequalities is decidable. This decidability result is quite striking, given the undecidability result proved by Matijacević [113] for a similar question, Hilbert’s 10-th problem: there is no algorithm deciding whether or not a general system of Diophantine equations has an integer solution.

In Chapter 3 we develop some elementary properties of semi-algebraic sets. Since we work over various real closed fields, and not only over the reals, it is necessary to reexamine several notions whose classical definitions break down in non-archimedean real closed fields. Examples of these are connectedness and compactness. Our proofs use non-archimedean real closed field extensions, which contain infinitesimal elements and can be described geometrically as germs of semi-algebraic functions, and algebraically as algebraic Puiseux series. The real closed field of algebraic Puiseux series plays a key role in the complexity results of Chapters 13 to 16.

Chapter 4 describes several algebraic results, relating in various ways properties of univariate and multivariate polynomials to linear algebra, determinants and quadratic forms. A general theme is to express some properties of univariate polynomials by the vanishing of specific polynomial expressions in their coefficients. The discriminant of a univariate polynomial  $P$ , for example, is a polynomial in the coefficients of  $P$  which vanishes when  $P$  has a multiple root. The discriminant is intimately related to real root counting, since, for polynomials of a fixed degree, all of whose roots are distinct, the sign of the discriminant determines the number of real roots modulo 4. The discriminant is in fact the determinant of a symmetric matrix whose signature gives an alternative method to Sturm’s for real root counting due to Hermite.

Similar polynomial expressions in the coefficients of two polynomials are the classical resultant and its generalization to subresultant coefficients. The vanishing of these subresultant coefficients expresses the fact that the greatest common divisor of two polynomials has at least a given degree. The resultant makes possible a constructive proof of a famous theorem of Hilbert, the Nullstellensatz, which provides a link between algebra and geometry in the algebraically closed case. Namely, the geometric statement ‘an algebraic variety (the common zeros of a finite family of polynomials) is empty’ is equivalent to the algebraic statement ‘1 belongs to the ideal generated by these polynomials’. An algebraic characterization of those systems of polynomial equations with a finite number of solutions in an algebraically closed field follows from Hilbert’s Nullstellensatz: a system of polynomial equations has a finite number of solutions in an algebraically closed field if and only if the corresponding quotient ring is a finite dimensional vector space. As seen in Chapter 1, the projection of an algebraic set in affine space is constructible. Considering projective space allows an even more satisfactory result: the projection of an algebraic set in projective space is algebraic. This result appears here as a consequence of a quantitative version of Hilbert’s Nullstellensatz, following the analysis of its constructive proof. A weak version of Bezout’s theorem, bounding the number of simple solutions of polynomials systems is a consequence of this projection theorem.

Semi-algebraic sets are defined by a finite number of polynomial inequalities. On the real line, semi-algebraic sets consist of a finite number of points and intervals. It is thus natural to wonder what kind of geometric finiteness properties are enjoyed by semi-algebraic sets in higher dimensions. In Chapter 5 we study various decompositions of a semi-algebraic set into a finite number of simple pieces. The most basic decomposition is called a cylindrical decomposition: a semi-algebraic set is decomposed into a finite number of pieces, each homeomorphic to an open cube. A finer decomposition provides a stratification, i.e. a decomposition into a finite number of pieces, called strata, which are smooth manifolds, such that the closure of a stratum is a union of strata of lower dimension. We also describe how to triangulate a closed and bounded semi-algebraic set. Various other finiteness results about semi-algebraic sets follow from these decompositions. Among these are:

- a semi-algebraic set has a finite number of connected components each of which is semi-algebraic,
- algebraic sets described by polynomials of fixed degree have a finite number of topological types.

A natural question raised by these results is to find explicit bounds on these quantities now known to be finite.

Chapter 6 is devoted to a self contained development of the basics of elementary algebraic topology. In particular, we define simplicial homology theory and, using the triangulation theorem, show how to associate to semi-algebraic sets certain discrete objects (the simplicial homology vector spaces) which are invariant under semi-algebraic homeomorphisms. The dimensions of these vector spaces, the Betti numbers, are an important measure of the topological complexity of semi-algebraic sets, the first of them being the number of connected components of the set. We also define the Euler-Poincaré characteristic, which is a significant topological invariant of algebraic and semi-algebraic sets.

Chapter 7 presents basic results of Morse theory and proves the classical Oleinik-Petrovsky-Thom-Milnor bounds on the sum of the Betti numbers of an algebraic set of a given degree. The basic technique for these results is the critical point method, which plays a key role in the complexity results of the last chapters of the book. According to basic results of Morse theory, the critical points of a well chosen projection on a line of a smooth hypersurface are precisely the places where a change in topology occurs in the part of the hypersurface inside a half space defined by a hyperplane orthogonal to the line. Counting these critical points using Bezout's theorem yields the Oleinik-Petrovsky-Thom-Milnor bound on the sum of the Betti numbers of an algebraic hypersurface, which is polynomial in the degree and exponential in the number of variables. More recent results bounding the individual Betti numbers of sign conditions defined by a family of polynomials on an algebraic set are described. These results involve a combinatorial part, depending on the number of polynomials considered, which is polynomial in the number of polynomials and exponential in the dimension of the algebraic set, and an algebraic part, given by the Oleinik-Petrovsky-Thom-Milnor bound. The combinatorial part of these bounds agrees with the number of connected components defined by a family of hyperplanes. These quantitative results on the number of connected components and Betti numbers of semi-algebraic sets provide an indication about the complexity results to be hoped for when studying various algorithmic problems related to semi-algebraic sets.

The second part of the book discusses various algorithmic problems in detail. These are mainly real root counting, deciding the existence of solutions for systems of equations and inequalities, computing the projection of a semi-algebraic set, deciding a sentence of the theory of real closed fields, eliminating quantifiers, and computing topological properties of algebraic and semi-algebraic sets.

In Chapter 8 we discuss a few notions of complexity needed to analyze our algorithms and discuss basic algorithms for linear algebra and remainder sequences. We perform a study of a useful tool closely related to remainder sequence, the subresultant sequence. This subresultant sequence plays an important role in modern methods for real root counting in Chapter 9, and

also provides a link between the classical methods of Sturm and Hermite seen earlier. Various methods for performing real root counting, and computing the signature of related quadratic forms, as well as an application to counting complex roots in a half plane, useful in control theory, are described.

Chapter 10 is devoted to real roots. In the field of the reals, which is archimedean, root isolation techniques are possible. They are based on Descartes's law of signs, presented in Chapter 2 and properties of Bernstein polynomials, which provide useful constructions in CAD (Computer Aided Design). For a general real closed field, isolation techniques are no longer possible. We prove that a root of a polynomial can be uniquely described by sign conditions on the derivatives of this polynomial, and we describe a different method for performing sign determination and characterizing real roots, without approximating the roots.

In Chapter 11, we describe an algorithm for computing the cylindrical decomposition which had been already studied in Chapter 5. The basic idea of this algorithm is to successively eliminate variables, using subresultants. Cylindrical decomposition has numerous applications among which are: deciding the truth of a sentence, eliminating quantifiers, computing a stratification, and computing topological information of various kinds, an example of which is computing the topology of an algebraic curve. The huge degree bounds (doubly exponential in the number of variables) output by the cylindrical decomposition method give estimates on the number of connected components of semi-algebraic sets which are much worse than those we obtained using the critical point method in Chapter 7.

The main idea developed in Chapters 12 to 16 is that, using the critical point method in an algorithmic way yields much better complexity bounds than those obtained by cylindrical decomposition for deciding the existential theory of the reals, eliminating quantifiers, deciding connectivity and computing connected components.

Chapter 12 is devoted to polynomial system solving. We give a few results about Gröbner bases, and explain the technique of rational univariate representation. Since our techniques in the following chapters involve infinitesimal deformations, we also indicate how to compute the limit of the bounded solutions of a polynomial system when the deformation parameters tend to zero. As a consequence, using the ideas of the critical point method described in Chapter 7, we are able to find a point in every connected components of an algebraic set. Since we deal with arbitrary algebraic sets which are not necessarily smooth, we introduce the notion of a pseudo-critical point in order to adapt the critical point method to this new situation. We compute a point in every semi-algebraically connected component of a bounded algebraic set with complexity polynomial in the degree and exponential in the number of variables. Using a similar technique, we compute the Euler-Poincaré characteristic of an algebraic set, with complexity polynomial in the degree and exponential in the number of variables.

In Chapter 13 we present an algorithm for the existential theory of the reals whose complexity is singly exponential in the number of variables. Using the pseudo-critical points introduced in Chapter 12 and perturbation methods to obtain polynomials in general position, we can compute the set of realizable sign conditions and compute representative points in each of the realizable sign conditions. Applications to the size of a ball meeting every connected component and various real and complex decision problems are provided. Finally we explain how to compute points in realizable sign conditions on an algebraic set taking advantage of the (possibly low) dimension of the algebraic set. We also compute the Euler-Poincaré characteristic of sign conditions defined by a set of polynomials. The complexity results obtained are quite satisfactory in view of the quantitative bounds proved in Chapter 7.

In Chapter 14 the results on the complexity of the general decision problem and quantifier elimination obtained in Chapter 11 using cylindrical decomposition are improved. The main idea is that the complexity of quantifier elimination should not be doubly exponential in the number of variables but rather in the number of blocks of variables appearing in the formula where the blocks of variables are delimited by alternations in the quantifiers  $\exists$  and  $\forall$ . The key notion is the set of realizable sign conditions of a family of polynomials for a given block structure of the set of variables, which is a generalization of the set of realizable sign conditions, corresponding to one single block. Parametrized versions of the methods presented in Chapter 13 give the technique needed for eliminating a whole block of variables.

In Chapters 15 and 16, we compute roadmaps and connected components of algebraic and semi-algebraic sets. Roadmaps can be intuitively described as an one dimensional skeleton of the set, providing a way to count connected components and to decide whether two points belong to the same connected component. A motivation for studying these problems comes from robot motion planning where the free space of a robot (the subspace of the configuration space of the robot consisting of those configurations where the robot is neither in conflict with its environment nor itself) can be modeled as a semi-algebraic set. In this context it is important to know whether a robot can move from one configuration to another. This is equivalent to deciding whether the two corresponding points in the free space are in the same connected component of the free space. The construction of roadmaps is based on the critical point method, using properties of pseudo-critical values. The complexity of the construction is singly exponential in the number of variables, which is a complexity much better than the one provided by cylindrical decomposition. Our construction of parametrized paths gives an algorithm for computing coverings of semi-algebraic sets by contractible sets, which in turn provides a single exponential time algorithm for computing the first Betti number of semi-algebraic sets. Moreover, it gives an efficient algorithm for computing semi-algebraic descriptions of the connected components of a semi-algebraic set in single exponential time.

*1 Warning* This book is intended to be self contained, assuming only that the reader has a basic knowledge of linear algebra and the rudiments of a basic course in algebra through the definitions and basic properties of groups, rings and fields, and in topology through the elementary properties of closed, open, compact and connected sets.

There are many other aspects of real algebraic geometry that are not considered in this book. The reader who wants to pursue the many aspects of real algebraic geometry beyond the introduction to the small part of it that we provide is encouraged to study other text books [26, 95, 5]. There is also a great deal of material about algorithms in real algebraic geometry that we are not covering in this book. To mention but a few: fewnomials, effective positivstellensatz, semi-definite programming, complexity of quadratic maps and quadratic sets, ...

*2 References* We have tried to keep our style as informal as possible. Rather than giving bibliographic references and footnotes in the body of the text, we have a section at the end of each chapter giving a brief description of the history of the results with a few of the relevant bibliographic citations. We only try to indicate where, to the best of our knowledge, the main ideas and results appear for the first time, and do not describe the full history and bibliography. We also list below the references containing the material we have used directly.

*3 Existing implementations* In terms of existing implementation of the algorithms described in the book, the current situation can be roughly summarized as follows: algorithms appearing in Chapters 8 to 12, or more efficient versions based on similar ideas, have been implemented (see a few references below). For most of the algorithms presented in Chapter 13 to 16, there is no implementation at all. The reason for that is that the methods developed are well adapted to complexity results but are not adapted to efficient implementation.

Most algorithms from Chapters 8 to 11 are quite classical and have been implemented several times. We refer to [40] since it is a recent implementation based directly on [20]. It uses in part the work presented in [29]. A very efficient variant of the real root isolation algorithm in the monomial basis in Chapter 10 is described in [138]. Cylindrical algebraic decomposition discussed in Chapter 11 has also been implemented many times, see for example [46, 30, 151]. We refer to [71] for an implementation of an algorithm computing the topology of real algebraic curves close to the one we present in Chapter 11. About algorithms discussed in Chapter 12, most computer algebra systems include Gröbner basis computations. Particularly efficient Gröbner basis computations, based on algorithms not described in the book, can be found in [59]. A very efficient rational univariate representation can be found in [135]. Computing a point in every connected component of an algebraic set based on critical point method techniques is done efficiently in [143], based on the algorithms developed in [8, 144].

*4 Comments about the second edition* An important change in content between the first edition [20] and the second one is the inversion of the order of Chapter 12 and Chapter 11. Indeed when teaching courses based on the book, we felt that the material on polynomial system solving was not necessary to explain cylindrical decomposition and it was better to make these two chapters independent for teaching purposes. For the same reason, we also made the real root counting technique based on signed subresultant coefficients independent of the signed subresultant polynomials and included it in Chapter 4 rather than in Chapter 9 as before. Some other chapters have been slightly reorganized. Several new topics are included in this second edition: results about normal polynomials and virtual roots in Chapter 2, about discriminants of symmetric matrices in Chapter 4, a new section bounding the Betti numbers of semi-algebraic sets in Chapter 7, an improved complexity analysis of real root isolation, as well as the real root isolation algorithm in the monomial basis, in Chapter 10, the notion of parametrized path in Chapter 15 and the computation of the first Betti number of a semi-algebraic set in single exponential time. We also included a table of notation and completed the bibliography and bibliographical notes at the end of the chapters. Various mistakes and typos have been corrected, and new ones introduced, for sure. As a result of the changes, the numbering of Definitions, Theorems etc. are not identical in the first edition [20] and the second one. Also, Algorithms now have their own numbering.

According to our contract with Springer-Verlag, we have had the right to post updated versions of the first edition of the book on our websites since December 2004. Currently an updated version of the first edition is available online as `bpr-posted1.pdf`. We are going to update on a regular basis this posted version. Here are the various url where these files can be obtained through `http://` at

`www.math.gatech.edu/~saugata/bpr-posted1.html`

`www.math.nyu.edu/faculty/pollack/bpr-posted1.html`

`perso.univ-rennes1.fr/marie-francoise.roy/bpr-posted1.html`

An implementation of algorithms from Chapters 8 to 10 and part of Chapter 11 written in Maxima by Fabrizio Caruso, as well as a version of Jean-Charles Faugère [59] and Fabrice Rouillier [135] software illustrating part of Chapter 12, can also be downloaded at `bpr-posted1-annex`.

Note that the second edition has been prepared inside  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$ . The  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$  files have been initially produced from classical latex files of the first edition. Even though some manual changes in the latex files have been necessary to obtain correct  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$  files, the translation into  $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$  was made automatically, and it has not been necessary to retype the text and formulas, besides a few exceptions.



After eighteen months of the publication of the current edition of the book, we will post the second edition online and it will be available for downloading from the same url as above.

*5 Interactive version of the book* Another possibility is to get the book as a  $\text{T}_{\text{E}}\text{X}_{\text{M}}\text{A}^{\text{C}}\text{S}$  project by downloading `bpr-posted1-int`. In the  $\text{T}_{\text{E}}\text{X}_{\text{M}}\text{A}^{\text{C}}\text{S}$  project version, you are able to travel in the book by clicking on references, to fold/unfold proofs, descriptions of the algorithms and parts of the text. You can use the open-source maxima code corresponding to algorithms of Chapters 8 to 10 and part of Chapter 11 written by Fabrizio Caruso [40]: check examples, read the source code and make your own computations inside the book. You can also use the part of [59] and [135] provided by Jean-Charles Faugère and Fabrice Rouillier to illustrate part of Chapter 12 directly in the book. These functionalities are still experimental. You are welcome to report to the authors' email addresses any problem you might meet in using them.

In the future,  $\text{T}_{\text{E}}\text{X}_{\text{M}}\text{A}^{\text{C}}\text{S}$  versions of the book will include other interactive features, such as being able to find all places in the book where a given theorem is quoted.

*6 Errors* If you find remaining errors in the book, we would appreciate it if you would let us know

email: `saugata.basu@math.gatech.edu`  
`pollack@cims.nyu.edu`  
`marie-francoise.roy@univ-rennes1.fr`

A list of errors identified in this version will be found at

`www.math.gatech.edu/~saugata/bpr_book/bpr-ed2-errata.html`.

*7 Acknowledgment* We thank Michel Coste, Greg Friedman, Laureano Gonzalez-Vega, Abdeljaoued Jounaidi, Henri Lombardi, Dimitri Pasechnik, Fabrice Rouillier for their advice and help. We also thank Solen Corvez, Gwenael Guérard, Michael Kettner, Tomas Lajous, Samuel Lelièvre, Mohab Safey, and Brad Weir for studying preliminary versions of the text and helping to improve it. Mistakes or typos in [20] have been identified by Morou Amidou, Emmanuel Briand, Fabrizio Caruso, Fernando Carreras, Keven Commault, Anne Devys, Arno Eigenwillig, Vincent Guenanff, Michael Kettner, Assia Mahboubi, Iona Necula, Adamou Otto, Dimitri Pasechnik, Hervé Perdry, Savvas Perikleous, Moussa Seydou.

Joris Van der Hoeven has provided support for the use of  $\text{T}_{\text{E}}\text{X}_{\text{M}}\text{A}^{\text{C}}\text{S}$  and produced several new versions of the software adapted to our purpose. Most figures are the same as in the first edition. However, Henri Lesourd produced some native  $\text{T}_{\text{E}}\text{X}_{\text{M}}\text{A}^{\text{C}}\text{S}$  diagrams and figures for us.

At different stages of writing this book the authors received support from CNRS, NSF, Université de Rennes 1, Courant Institute of Mathematical Sciences, University of Michigan, Georgia Institute of Technology, the RIP Program in Oberwolfach, MSRI, DIMACS, RISC, Linz, Centre Emile Borel. Fabrizio Caruso was supported by RAAG during a post doctoral fellowship in Rennes and Santander. The software due to Jean-Charles Faugère [59] and Fabrice Rouillier [135] was developed under the SALSA project at INRIA, CNRS and Université Pierre et Marie Curie, Paris.

*8 Sources* Our sources for Chapter 2 are: [26] for Section 2.1 and Section 2.4, [140, 98, 49] for Section 2.2, [47] for Section 2.3 and [164, 109] for Section 2.5. Our source for Section 3.1, Section 3.2 and Section 3.3 of Chapter 3 is [26]. Our sources for Chapter 4 are: [63] for Section 4.1, [94] for Theorem 4.47 in Section 4.4, [159, 147] for Section 4.4, [128, 129] for Section 4.6 and [22] for Section 4.7. Our sources for Chapter 5 are [26, 47, 48]. Our source for Chapter 6 is [150]. Our sources for Chapter 7 are [117, 26, 17], and for Section 7.5 [62, 21]. Our sources for Chapter 8 are: [1] for Section 8.2 and [112] for Section 8.3. Our sources for Chapter 9 are [63] and [66, 69, 70, 140, 2] for part of Section 9.1. Our sources for Chapter 10 are: [116] for Section 10.1, [138, 149] for Section 10.2, [141] for Sections 10.3 and [129] for Section 10.4. Our source for Section 11.4 is [52], and for Section 11.6 is [67]. Our sources for Chapter 12 are: for Section 12.1 [51], for Section 12.2 [72], for Section 12.4 [4, 134], for Section 12.5 [13]. The results presented in Section 13.1, Section 13.2 and Section 13.3 of Chapter 13 are based on [13, 15]. Our source for Section 13.4 of Chapter 13 is [18]. Our source for Chapter 14 is [13]. Our sources for Chapter 15 and Chapter 16 are [16, 21].

---

## Algebraically Closed Fields

The main purpose of this chapter is the definition of constructible sets and the statement that, in the context of algebraically closed fields, the projection of a constructible set is constructible.

Section 1.1 is devoted to definitions. The main technique used for proving the projection theorem in Section 1.3 is the remainder sequence defined in Section 1.2 and, for the case where the coefficients have parameters, the tree of possible pseudo-remainder sequences. Several important applications of logical nature of the projection theorem are given in Section 1.4.

### 1.1 Definitions and First Properties

The objects of our interest in this section are sets defined by polynomials with coefficients in an algebraically closed field  $C$ .

A field  $C$  is **algebraically closed** if any non-constant univariate polynomial  $P(X)$  with coefficients in  $C$  has a **root** in  $C$ , i.e. there exists  $x \in C$  such that  $P(x) = 0$ .

Every field has a minimal extension which is algebraically closed and this extension is called the **algebraic closure** of the field (see Section 2, Chapter 5 of [102]). A typical example of an algebraically closed field is the field  $\mathbb{C}$  of complex numbers.

We study the sets of points which are the common zeros of a finite family of polynomials.

If  $D$  is a ring, we denote by  $D[X_1, \dots, X_k]$  the polynomials in  $k$  variables  $X_1, \dots, X_k$  with coefficients in  $D$ .

**Notation 1.1. [Zero set]** If  $\mathcal{P}$  is a finite subset of  $C[X_1, \dots, X_k]$  we write the **set of zeros of  $\mathcal{P}$  in  $C^k$**  as

$$\text{Zer}(\mathcal{P}, C^k) = \{x \in C^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\}.$$

These are the **algebraic subsets** of  $C^k$ .

The set  $C^k$  is algebraic since  $C^k = \text{Zer}(\{0\}, C^k)$ . □

**Exercise 1.1.** Prove that an algebraic subset of  $C$  is either a finite set or empty or equal to  $C$ .

It is natural to consider the smallest family of sets which contain the algebraic sets and is also closed under the boolean operations (complementation, finite unions, and finite intersections). These are the **constructible sets**. Similarly, the smallest family of sets which contain the algebraic sets, their complements, and is closed under finite intersections is the family of **basic constructible sets**. Such a basic constructible set  $S$  can be described as a conjunction of polynomial equations and inequations, namely

$$S = \{x \in C^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(x) \neq 0\}$$

with  $\mathcal{P}, \mathcal{Q}$  finite subsets of  $C[X_1, \dots, X_k]$ .

**Exercise 1.2.** Prove that a constructible subset of  $C$  is either a finite set or the complement of a finite set.

**Exercise 1.3.** Prove that a constructible set in  $C^k$  is a finite union of basic constructible sets.

The principal goal of this chapter is to prove that the projection from  $C^{k+1}$  to  $C^k$  that is defined by "forgetting" the last coordinate maps constructible sets to constructible sets. For this, since projection commutes with union, it suffices to prove that the projection

$$\{y \in C^k \mid \exists x \in C \bigwedge_{P \in \mathcal{P}} P(y, x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(y, x) \neq 0\}$$

of a basic constructible set,

$$\{(y, x) \in C^{k+1} \mid \bigwedge_{P \in \mathcal{P}} P(y, x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(y, x) \neq 0\}$$

is constructible, i.e. can be described by a boolean combination of polynomial **equations** ( $P = 0$ ) and **inequations** ( $P \neq 0$ ) in  $Y = (Y_1, \dots, Y_k)$ .

Some terminology from logic is useful for the study of constructible sets.

We define the language of fields by describing the formulas of this language. The formulas are built starting with atoms, which are polynomial equations and inequations. A formula is written using atoms together with the logical connectives "and", "or", and "negation" ( $\wedge$ ,  $\vee$ , and  $\neg$ ) and the existential and universal quantifiers ( $\exists$ ,  $\forall$ ). A formula has free variables, i.e. non-quantified variables, and bound variables, i.e. quantified variables. More precisely, let  $D$  be a subring of  $C$ . We define the **language of fields with coefficients in  $D$**  as follows. An **atom** is  $P = 0$  or  $P \neq 0$ , where  $P$  is a polynomial in  $D[X_1, \dots, X_k]$ . We define simultaneously the **formulas** and  $\text{Free}(\Phi)$ , the set of **free variables of a formula  $\Phi$** , as follows

- an atom  $P = 0$  or  $P \neq 0$ , where  $P$  is a polynomial in  $D[X_1, \dots, X_k]$  is a formula with free variables  $\{X_1, \dots, X_k\}$ ,

- if  $\Phi_1$  and  $\Phi_2$  are formulas, then  $\Phi_1 \wedge \Phi_2$  and  $\Phi_1 \vee \Phi_2$  are formulas with

$$\text{Free}(\Phi_1 \wedge \Phi_2) = \text{Free}(\Phi_1 \vee \Phi_2) = \text{Free}(\Phi_1) \cup \text{Free}(\Phi_2),$$

- if  $\Phi$  is a formula, then  $\neg(\Phi)$  is a formula with

$$\text{Free}(\neg(\Phi)) = \text{Free}(\Phi),$$

- if  $\Phi$  is a formula and  $X \in \text{Free}(\Phi)$ , then  $(\exists X) \Phi$  and  $(\forall X) \Phi$  are formulas with

$$\text{Free}((\exists X) \Phi) = \text{Free}((\forall X) \Phi) = \text{Free}(\Phi) \setminus \{X\}.$$

If  $\Phi$  and  $\Psi$  are formulas,  $\Phi \Rightarrow \Psi$  is the formula  $\neg(\Phi) \vee \Psi$ .

A **quantifier free formula** is a formula in which no quantifier appears, neither  $\exists$  nor  $\forall$ . A **basic formula** is a conjunction of atoms.

The **C-realization of a formula**  $\Phi$  with free variables contained in  $\{Y_1, \dots, Y_k\}$ , denoted  $\text{Reali}(\Phi, C^k)$ , is the set of  $y \in C^k$  such that  $\Phi(y)$  is true. It is defined by induction on the construction of the formula, starting from atoms:

$$\begin{aligned} \text{Reali}(P = 0, C^k) &= \{y \in C^k \mid P(y) = 0\}, \\ \text{Reali}(P \neq 0, C^k) &= \{y \in C^k \mid P(y) \neq 0\}, \\ \text{Reali}(\Phi_1 \wedge \Phi_2, C^k) &= \text{Reali}(\Phi_1, C^k) \cap \text{Reali}(\Phi_2, C^k), \\ \text{Reali}(\Phi_1 \vee \Phi_2, C^k) &= \text{Reali}(\Phi_1, C^k) \cup \text{Reali}(\Phi_2, C^k), \\ \text{Reali}(\neg\Phi, C^k) &= C^k \setminus \text{Reali}(\Phi, C^k), \\ \text{Reali}((\exists X) \Phi, C^k) &= \{y \in C^k \mid \exists x \in C \quad (x, y) \in \text{Reali}(\Phi, C^{k+1})\}, \\ \text{Reali}((\forall X) \Phi, C^k) &= \{y \in C^k \mid \forall x \in C \quad (x, y) \in \text{Reali}(\Phi, C^{k+1})\} \end{aligned}$$

Two formulas  $\Phi$  and  $\Psi$  such that  $\text{Free}(\Phi) = \text{Free}(\Psi) = \{Y_1, \dots, Y_k\}$  are **C-equivalent** if  $\text{Reali}(\Phi, C^k) = \text{Reali}(\Psi, C^k)$ .

If there is no ambiguity, we simply write  $\text{Reali}(\Phi)$  for  $\text{Reali}(\Phi, C^k)$  and talk about realization and equivalence.

*Example 1.2.* The formulas  $\Phi = ((\exists Y) XY - 1 = 0)$  and  $\Psi = (X \neq 0)$  are two formulas of the language of fields with coefficients in  $\mathbb{Z}$  and

$$\text{Free}(\Phi) = \text{Free}(\Psi) = \{X\}.$$

Note that the formula  $\Psi$  is quantifier free. Moreover,  $\Phi$  and  $\Psi$  are C-equivalent since

$$\begin{aligned} \text{Reali}(\Phi, C) &= \{x \in C \mid \exists y \in C \quad xy - 1 = 0\} \\ &= \{x \in C \mid x \neq 0\} \\ &= \text{Reali}(\Psi, C). \end{aligned}$$

□

It is clear that a set is constructible if and only if it can be represented as the realization of a quantifier free formula.

It is easy to see that any formula  $\Phi$  with  $\text{Free}(\Phi) = \{Y_1, \dots, Y_k\}$  in the language of fields with coefficients in  $D$  is  $C$ -equivalent to a formula

$$(\text{Qu}_1 X_1) \dots (\text{Qu}_m X_m) \mathcal{B}(X_1, \dots, X_m, Y_1, \dots, Y_k)$$

where each  $\text{Qu}_i \in \{\forall, \exists\}$  and  $\mathcal{B}$  is a quantifier free formula involving polynomials in  $D[X_1, \dots, X_m, Y_1, \dots, Y_k]$ . This is called its **prenex normal form** (see Section 10, Chapter 1 of [115]). The variables  $X_1, \dots, X_m$  are called **bound variables**.

If the formula  $\Phi$  has no free variables, i.e.  $\text{Free}(\Phi) = \emptyset$ , then it is called a **sentence**, and it is either  $C$ -equivalent to true, when  $\text{Reali}(\Phi), \{0\} = \{0\}$ , or  $C$ -equivalent to false, when  $\text{Reali}(\Phi), \{0\} = \emptyset$ . For example,  $0 = 0$  is  $C$ -equivalent to true, and  $0 = 1$  is  $C$ -equivalent to false.

*Remark 1.3.* Though many statements of algebra can be expressed by a sentence in the language of fields, it is necessary to be careful in the use of this notion. Consider for example the fundamental theorem of algebra: any non constant polynomial with coefficients in  $\mathbb{C}$  has a root in  $\mathbb{C}$ , which is expressed by

$$\forall P \in \mathbb{C}[X] \text{ deg}(P) > 0, \exists X \in \mathbb{C} P(X) = 0.$$

This expression is not a sentence of the language of fields with coefficients in  $\mathbb{C}$ , since quantification over all polynomials is not allowed in the definition of formulas. However, fixing the degree to be equal to  $d$ , it is possible to express by a sentence  $\Phi_d$  the statement: any monic polynomial of degree  $d$  with coefficients in  $\mathbb{C}$  has a root in  $\mathbb{C}$ . We write as an example

$$\Phi_2 = ((\forall Y_1) (\forall Y_2) (\exists X) X^2 + Y_1 X + Y_2 = 0).$$

So the definition of an algebraically closed field can be expressed by an infinite list of sentences in the language of fields: the field axioms and the sentences  $\Phi_d, d \geq 1$ .  $\square$

**Exercise 1.4.** Write the formulas for the axioms of fields.

## 1.2 Euclidean Division and Greatest Common Divisor

We study euclidean division, compute greatest common divisors, and show how to use them to decide whether or not a basic constructible set of  $C$  is empty.

In this section,  $C$  is an algebraically closed field,  $D$  a subring of  $C$  and  $K$  the quotient field of  $D$ . One can take as a typical example of this situation the field  $\mathbb{C}$  of complex numbers, the ring  $\mathbb{Z}$  of integers, and the field  $\mathbb{Q}$  of rational numbers.

Let  $P$  be a non-zero polynomial

$$P = a_p X^p + \cdots + a_1 X + a_0 \in D[X]$$

with  $a_p \neq 0$ .

We denote the **degree of  $P$** , which is  $p$ , by  $\deg(P)$ . By convention, the degree of the zero polynomial is defined to be  $-\infty$ . If  $P$  is non-zero, we write  $\text{coef}_j(P) = a_j$  for the **coefficient of  $X^j$  in  $P$**  (which is equal to 0 if  $j > \deg(P)$ ) and  $\text{lcoef}(P)$  for its **leading coefficient**  $a_p = \text{coef}_{\deg(P)}(P)$ . By convention  $\text{lcoef}(0) = 1$ .

Suppose that  $P$  and  $Q$  are two polynomials in  $D[X]$ . The polynomial  $Q$  is a **divisor of  $P$**  if  $P = AQ$  for some  $A \in K[X]$ . Thus, while every  $P$  divides 0, 0 divides 0 and no other polynomial.

If  $Q \neq 0$ , the **remainder** in the **euclidean division of  $P$  by  $Q$** , denoted  $\text{Rem}(P, Q)$ , is the unique polynomial  $R \in K[X]$  of degree smaller than the degree of  $Q$  such that  $P = AQ + R$  with  $A \in K[X]$ . The **quotient** in the euclidean division of  $P$  by  $Q$ , denoted  $\text{Quo}(P, Q)$ , is  $A$ .

**Exercise 1.5.** Prove that, if  $Q \neq 0$ , there exists a unique pair  $(R, A)$  of polynomials in  $K[X]$  such that  $P = AQ + R$ ,  $\deg(R) < \deg(Q)$ .

*Remark 1.4.* Clearly,  $\text{Rem}(aP, bQ) = a\text{Rem}(P, Q)$  for any  $a, b \in K$  with  $b \neq 0$ . At a root  $x$  of  $Q$ ,  $\text{Rem}(P, Q)(x) = P(x)$ .  $\square$

**Exercise 1.6.** Prove that  $x$  is a root of  $P$  in  $K$  if and only if  $X - x$  is a divisor of  $P$  in  $K[X]$ .

**Exercise 1.7.** Prove that if  $C$  is algebraically closed, every  $P \in C[X]$  can be written uniquely as

$$P = a(X - x_1)^{\mu_1} \cdots (X - x_k)^{\mu_k},$$

with  $x_1, \dots, x_k$  distinct elements of  $C$ .

A **greatest common divisor of  $P$  and  $Q$** , denoted  $\text{gcd}(P, Q)$ , is a polynomial  $G \in K[X]$  such that  $G$  is a divisor of both  $P$  and  $Q$ , and any divisor of both  $P$  and  $Q$  is a divisor of  $G$ . Observe that this definition implies that  $P$  is a greatest common divisor of  $P$  and 0. Clearly, any two greatest common divisors (say  $G_1, G_2$ ) of  $P$  and  $Q$  must divide each other and have equal degree. Hence  $G_1 = aG_2$  for some  $a \in K$ . Thus, any two greatest common divisors of  $P$  and  $Q$  are proportional by an element in  $K \setminus \{0\}$ . Two polynomials are **coprime** if their greatest common divisor is an element of  $K \setminus \{0\}$ .

A **least common multiple of  $P$  and  $Q$** ,  $\text{lcm}(P, Q)$  is a polynomial  $G \in K[X]$  such that  $G$  is a multiple of both  $P$  and  $Q$ , and any multiple of both  $P$  and  $Q$  is a multiple of  $G$ . Clearly, any two least common multiples  $L_1, L_2$  of  $P$  and  $Q$  must divide each other and have equal degree. Hence  $L_1 = aL_2$  for some  $a \in K$ . Thus, any two least common multiple of  $P$  and  $Q$  are proportional by an element in  $K \setminus \{0\}$ .

It follows immediately from the definitions that:

**Proposition 1.5.** *Let  $P \in K[X]$  and  $Q \in K[X]$ , not both zero. Then  $PQ/G$  is a least common multiple of  $P$  and  $Q$ .*

**Corollary 1.6.**

$$\deg(\text{lcm}(P, Q)) = \deg(P) + \deg(Q) - \deg(\text{gcd}(P, Q)).$$

We now prove that greatest common divisors and least common multiple exist by using euclidean division repeatedly.

**Definition 1.7. [Signed remainder sequence]** Given  $P, Q \in K[X]$ , not both 0, we define the **signed remainder sequence of  $P$  and  $Q$** ,

$$\text{SRemS}(P, Q) = \text{SRemS}_0(P, Q), \text{SRemS}_1(P, Q), \dots, \text{SRemS}_k(P, Q)$$

by

$$\begin{aligned} \text{SRemS}_0(P, Q) &= P, \\ \text{SRemS}_1(P, Q) &= Q, \\ \text{SRemS}_2(P, Q) &= -\text{Rem}(\text{SRemS}_0(P, Q), \text{SRemS}_1(P, Q)), \\ &\vdots \\ \text{SRemS}_k(P, Q) &= -\text{Rem}(\text{SRemS}_{k-2}(P, Q), \text{SRemS}_{k-1}(P, Q)) \neq 0, \\ \text{SRemS}_{k+1}(P, Q) &= -\text{Rem}(\text{SRemS}_{k-1}(P, Q), \text{SRemS}_k(P, Q)) = 0. \end{aligned}$$

The signs introduced here are unimportant in the algebraically closed case. They play an important role when we consider analogous problems over real closed fields in Chapter 2.  $\square$

In the above, each  $\text{SRemS}_i(P, Q)$  is the negative of the remainder in the euclidean division of  $\text{SRemS}_{i-2}(P, Q)$  by  $\text{SRemS}_{i-1}(P, Q)$  for  $2 \leq i \leq k+1$ , and the sequence ends with  $\text{SRemS}_k(P, Q)$  when  $\text{SRemS}_{k+1}(P, Q) = 0$ , for  $k \geq 0$ .

**Proposition 1.8.** *The polynomial  $\text{SRemS}_k(P, Q)$  is a greatest common divisor of  $P$  and  $Q$ .*

**Proof:** Observe that if a polynomial  $A$  divides two polynomials  $B, C$  then it also divides  $UB + VC$  for arbitrary polynomials  $U, V$ . Since

$$\text{SRemS}_{k+1}(P, Q) = -\text{Rem}(\text{SRemS}_{k-1}(P, Q), \text{SRemS}_k(P, Q)) = 0,$$

$\text{SRemS}_k(P, Q)$  divides  $\text{SRemS}_{k-1}(P, Q)$  and since,

$$\text{SRemS}_{k-2}(P, Q) = -\text{SRemS}_k(P, Q) + A \text{SRemS}_{k-1}(P, Q),$$

$\text{SRemS}_k(P, Q)$  divides  $\text{SRemS}_{k-2}(P, Q)$  using the above observation. Continuing this process one obtains that  $\text{SRemS}_k(P, Q)$  divides  $\text{SRemS}_1(P, Q) = Q$  and  $\text{SRemS}_0(P, Q) = P$ .



Also, if any polynomial divides  $\text{SRemS}_0(P, Q)$ ,  $\text{SRemS}_1(P, Q)$  (that is  $P, Q$ ) then it divides  $\text{SRemS}_2(P, Q)$  and hence  $\text{SRemS}_3(P, Q)$  and so on. Hence, it divides  $\text{SRemS}_k(P, Q)$ .  $\square$

Note that the signed remainder sequence of  $P$  and  $0$  is  $P$  and when  $Q$  is not  $0$ , the signed remainder sequence of  $0$  and  $Q$  is  $0, Q$ .

Also, note that by unwinding the definitions of the  $\text{SRemS}_i(P, Q)$ , we can express  $\text{SRemS}_k(P, Q) = \gcd(P, Q)$  as  $UP + VQ$  for some polynomials  $U, V$  in  $K[X]$ . We prove bounds on the degrees of  $U, V$  by elucidating the preceding remark.

**Proposition 1.9.** *If  $G$  is a greatest common divisor of  $P$  and  $Q$ , then there exist  $U$  and  $V$  with*

$$UP + VQ = G.$$

Moreover, if  $\deg(G) = g$ ,  $U$  and  $V$  can be chosen so that  $\deg(U) < q - g$ ,  $\deg(V) < p - g$ .

The proof uses the extended signed remainder sequence defined as follows.

**Definition 1.10. [Extended signed remainder sequence]**

Given  $P, Q \in K[X]$ , not both  $0$ , let

$$\begin{aligned} \text{SRemU}_0(P, Q) &= 1, \\ \text{SRemV}_0(P, Q) &= 0, \\ \text{SRemU}_1(P, Q) &= 0, \\ \text{SRemV}_1(P, Q) &= 1, \\ A_{i+1} &= \text{Quo}(\text{SRemS}_{i-1}(P, Q), \text{SRemS}_i(P, Q)), \\ \text{SRemS}_{i+1}(P, Q) &= -\text{SRemS}_{i-1}(P, Q) + A_{i+1} \text{SRemS}_i(P, Q), \\ \text{SRemU}_{i+1}(P, Q) &= -\text{SRemU}_{i-1}(P, Q) + A_{i+1} \text{SRemU}_i(P, Q), \\ \text{SRemV}_{i+1}(P, Q) &= -\text{SRemV}_{i-1}(P, Q) + A_{i+1} \text{SRemV}_i(P, Q) \end{aligned}$$

for  $0 \leq i \leq k$  where  $k$  is the least non-negative integer such that  $\text{SRemS}_{k+1} = 0$ .

The **extended signed remainder sequence**  $\text{Ex}(P, Q)$  of  $P$  and  $Q$  is  $\text{Ex}_0(P, Q), \dots, \text{Ex}_k(P, Q)$  with

$$\text{Ex}_i(P, Q) = (\text{SRemS}_i(P, Q), \text{SRemU}_i(P, Q), \text{SRemV}_i(P, Q)). \quad \square$$

The proof of Proposition 1.9 uses the following lemma.

**Lemma 1.11.** *For  $0 \leq i \leq k + 1$ ,*

$$\text{SRemS}_i(P, Q) = \text{SRemU}_i(P, Q)P + \text{SRemV}_i(P, Q)Q.$$

Let  $d_i = \deg(\text{SRemS}_i(P, Q))$ . For  $1 \leq i \leq k$ ,  $\deg(\text{SRemU}_{i+1}(P, Q)) = q - d_i$ , and  $\deg(\text{SRemV}_{i+1}(P, Q)) = p - d_i$ .

**Proof:** It is easy to verify by induction on  $i$  that, for  $0 \leq i \leq k+1$ ,

$$\text{SRemS}_i(P, Q) = \text{SRemU}_i(P, Q)P + \text{SRemV}_i(P, Q)Q.$$

Note that  $d_i < d_{i-1}$ . The proof of the claim on the degrees proceeds by induction. Clearly, since

$$\begin{aligned} \text{SRemU}_2(P, Q) &= -1 \\ \text{SRemU}_3(P, Q) &= -\text{Quo}(\text{SRemS}_1(P, Q), \text{SRemS}_2(P, Q)), \\ \deg(\text{SRemU}_2(P, Q)) &= q - d_1, \\ \deg(\text{SRemU}_3(P, Q)) &= q - d_2. \end{aligned}$$

Similarly,

$$\begin{aligned} \deg(\text{SRemV}_2(P, Q)) &= p - d_1, \\ \deg(\text{SRemV}_3(P, Q)) &= p - d_2. \end{aligned}$$

Using the definitions of  $\text{SRemU}_{i+1}(P, Q)$ ,  $\text{SRemV}_{i+1}(P, Q)$  and the induction hypothesis, we get

$$\begin{aligned} \deg(\text{SRemU}_{i-1}(P, Q)) &= q - d_{i-2}, \\ \deg(\text{SRemU}_i(P, Q)) &= q - d_{i-1} \\ \deg(A_{i+1} \text{SRemU}_i(P, Q)) &= d_{i-1} - d_i + q - d_{i-1} \\ &= q - d_i > q - d_{i-2}. \end{aligned}$$

Hence,  $\deg(\text{SRemU}_{i+1}) = q - d_i$ . Similarly,

$$\begin{aligned} \deg(\text{SRemV}_{i-1}(P, Q)) &= p - d_{i-2}, \\ \deg(\text{SRemV}_i(P, Q)) &= p - d_{i-1} \\ \deg(A_{i+1} \text{SRemV}_i(P, Q)) &= d_{i-1} - d_i + p - d_{i-1} \\ &= p - d_i > p - d_{i-2}. \end{aligned}$$

Hence,  $\deg(\text{SRemV}_{i+1}(P, Q)) = p - d_i$ . □

**Proof of Proposition 1.9:** The claim follows by Lemma 1.11 and Proposition 1.8 since  $\text{SRemS}_k(P, Q)$  is a gcd of  $P$  and  $Q$ , taking

$$U = \text{SRemU}_k(P, Q), V = \text{SRemV}_k(P, Q),$$

and noting that  $p - d_{k-1} < p - g$ ,  $q - d_{k-1} < q - g$ . □

The extended signed remainder sequence also provides a least common multiple of  $P$  and  $Q$ .

**Proposition 1.12.** *The equality*

$$\text{SRemU}_{k+1}(P, Q)P = -\text{SRemV}_{k+1}(P, Q)Q.$$

*holds and  $\text{SRemU}_{k+1}(P, Q)P = -\text{SRemV}_{k+1}(P, Q)Q$  is a least common multiple of  $P$  and  $Q$ .*

**Proof:** Since  $d_k = \deg(\gcd(P, Q))$ ,  $\deg(\text{SRem}U_{k+1}(P, Q)) = q - d_k$ ,  $\deg(\text{SRem}V_k(P, Q)) = p - d_k$ , and

$$\text{SRem}U_{k+1}(P, Q)P + \text{SRem}V_{k+1}(P, Q)Q = 0,$$

it follows that

$$\text{SRem}U_{k+1}(P, Q)P = -\text{SRem}V_{k+1}(P, Q)Q$$

is a common multiple of  $P$  and  $Q$  of degree  $p + q - d_k$ , hence a least common multiple of  $P$  and  $Q$ . □

**Definition 1.13. [Greatest common divisor of a family]** A **greatest common divisor of a finite family of polynomials** is a divisor of all the polynomials in the family that is also a multiple of any polynomial that divides every polynomial in the family. A greatest common divisor of a family can be obtained inductively on the number of elements of the family by

$$\begin{aligned} \gcd(\emptyset) &= 0, \\ \gcd(\mathcal{P} \cup \{P\}) &= \gcd(P, \gcd(\mathcal{P})). \end{aligned}$$

□

Note that

- $x \in C$  is a root of every polynomial in  $\mathcal{P}$  if and only if it is a root of  $\gcd(\mathcal{P})$ ,
- $x \in C$  is not a root of any polynomial in  $\mathcal{Q}$  if and only if it is not a root of  $\prod_{Q \in \mathcal{Q}} Q$  (with the convention that the product of the empty family is 1),
- every root of  $P$  in  $C$  is a root of  $Q$  if and only if  $\gcd(P, Q^{\deg(P)}) = P$  (with the convention that  $Q^{\deg(0)} = 0$ ).

With these observations the following lemma is clear:

**Lemma 1.14.** *If  $\mathcal{P}, \mathcal{Q}$  are two finite subsets of  $D[X]$ , then there is an  $x \in C$  such that*

$$\left( \bigwedge_{P \in \mathcal{P}} P(x) = 0 \right) \wedge \left( \bigwedge_{Q \in \mathcal{Q}} Q(x) \neq 0 \right)$$

*if and only if*

$$\deg(\gcd(\gcd(\mathcal{P}), \prod_{Q \in \mathcal{Q}} Q^d)) \neq \deg(\gcd(\mathcal{P})),$$

*where  $d$  is any integer greater than  $\deg(\gcd(\mathcal{P}))$ .*

Note that when  $\mathcal{Q} = \emptyset$ , since  $\prod_{Q \in \emptyset} Q = 1$ , the lemma says that there is an  $x \in C$  such that  $\bigwedge_{P \in \mathcal{P}} P(x) = 0$  if and only if  $\deg(\gcd(\mathcal{P})) \neq 0$ . Note also that when  $\mathcal{P} = \emptyset$ , the lemma says that there is an  $x \in C$  such that  $\bigwedge_{Q \in \mathcal{Q}} Q(x) \neq 0$  if and only if  $\deg(\prod_{Q \in \mathcal{Q}} Q) \geq 0$ , i.e.  $1 \notin \mathcal{Q}$ .

**Exercise 1.8.** Design an algorithm to decide whether or not a basic constructible set in  $C$  is empty.

### 1.3 Projection Theorem for Constructible Sets

Now that we know how to decide whether or not a basic constructible set in  $\mathbb{C}$  is empty, we can show that the projection from  $\mathbb{C}^{k+1}$  to  $\mathbb{C}^k$  of a basic constructible set is constructible. We shall do this by viewing the multivariate situation as a univariate situation with parameters. Viewing a univariate algorithm parametrically to obtain a multivariate algorithm is among the most important paradigms used throughout this book.

More precisely, the basic constructible set  $S \subset \mathbb{C}^{k+1}$  can be described as

$$S = \{z \in \mathbb{C}^{k+1} \mid \bigwedge_{P \in \mathcal{P}} P(z) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(z) \neq 0\}$$

with  $\mathcal{P}, \mathcal{Q}$  finite subsets of  $\mathbb{C}[Y_1, \dots, Y_k, X]$ , and its projection  $\pi(S)$  (forgetting the last coordinate) is

$$\pi(S) = \{y \in \mathbb{C}^k \mid \exists x \in \mathbb{C} \left( \bigwedge_{P \in \mathcal{P}} P(y, x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(y, x) \neq 0 \right)\}.$$

We can consider the polynomials in  $\mathcal{P}$  and  $\mathcal{Q}$  as polynomials in the single variable  $X$  with the variables  $(Y_1, \dots, Y_k)$  appearing as parameters. For a specialization of  $Y$  to  $y = (y_1, \dots, y_k) \in \mathbb{C}^k$ , we write  $P_y(X)$  for  $P(y_1, \dots, y_k, X)$ . Hence,

$$\pi(S) = \{y \in \mathbb{C}^k \mid \exists x \in \mathbb{C} \left( \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q_y(x) \neq 0 \right)\},$$

and, for a particular  $y \in \mathbb{C}^k$  we can decide, using Exercise 1.8, whether or not

$$\exists x \in \mathbb{C} \left( \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q_y(x) \neq 0 \right)$$

is true.

Defining

$$S_y = \{x \in \mathbb{C} \mid \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q_y(x) \neq 0\},$$

what is crucial now is to partition the parameter space  $\mathbb{C}^k$  into finitely many parts so that the decision algorithm testing whether  $S_y$  is empty or not is the same (is uniform) for all  $y$  in any given part. Because of this uniformity, it will turn out that each part of the partition is a constructible set. Since  $\pi(S)$  is the union of those parts where  $S_y \neq \emptyset$ ,  $\pi(S)$  is constructible being the union of finitely many constructible sets.

We next study the signed remainder sequence of  $P_y$  and  $Q_y$  for all possible specialization of  $Y$  to  $y \in \mathbb{C}^k$ . This cannot be done in a completely uniform way, since denominators appear in the euclidean division process. Nevertheless, fixing the degrees of the polynomials in the signed remainder sequence, it is possible to partition the parameter space,  $\mathbb{C}^k$ , into a finite number of parts so that the signed remainder sequence is uniform in each part.

*Example 1.15.* We consider a general polynomial of degree 4. Dividing by its leading coefficient, it is not a loss of generality to take  $P$  to be monic. So let  $P = X^4 + \alpha X^3 + \beta X^2 + \gamma X + \delta$ . Since the translation  $X \mapsto X - \alpha/4$  kills the term of degree 3, we can suppose  $P = X^4 + aX^2 + bX + c$ .

Consider  $P = X^4 + aX^2 + bX + c$  and its derivative  $P' = 4X^3 + 2aX + b$ . Their signed remainder sequence in  $\mathbb{Q}(a, b, c)[X]$  is

$$\begin{aligned} P &= X^4 + aX^2 + bX + c \\ P' &= 4X^3 + 2aX + b \\ S_2 &= -\text{Rem}(P, P') \\ &= -\frac{1}{2}aX^2 - \frac{3}{4}bX - c \\ S_3 &= -\text{Rem}(P', S_2) \\ &= \frac{(8ac - 9b^2 - 2a^3)X}{a^2} - \frac{b(12c + a^2)}{a^2} \\ S_4 &= -\text{Rem}(S_2, S_3) \\ &= \frac{1}{4} \frac{a^2(256c^3 - 128a^2c^2 + 144acb^2 - 16a^4c - 27b^4 - 4b^2a^3)}{(8ac - 9b^2 - 2a^3)^2} \end{aligned}$$

Note that when  $(a, b, c)$  are specialized to values in  $\mathbb{C}^3$  for which  $a = 0$  or  $8ac - 9b^2 - 2a^3 = 0$ , the signed remainder sequence of  $P$  and  $P'$  for these special values is not obtained by specializing  $a, b, c$  in the signed remainder sequence in  $\mathbb{Q}(a, b, c)[X]$ .  $\square$

In order to take into account all the possible signed remainder sequences that can appear when we specialize the parameters, we introduce the following definitions and notation.

We get rid of denominators appearing in the remainders through the notion of signed pseudo-remainders. Let

$$\begin{aligned} P &= a_p X^p + \cdots + a_0 \in D[X], \\ Q &= b_q X^q + \cdots + b_0 \in D[X], \end{aligned}$$

where  $D$  is a subring of  $C$ . Note that the only denominators occurring in the euclidean division of  $P$  by  $Q$  are  $b_q^i$ ,  $i \leq p - q + 1$ . The **signed pseudo-remainder** denoted  $\text{PRem}(P, Q)$ , is the remainder in the euclidean division of  $b_q^d P$  by  $Q$ , where  $d$  is the smallest even integer greater than or equal to  $p - q + 1$ . Note that the euclidean division of  $b_q^d P$  by  $Q$  can be performed in  $D$  and that  $\text{PRem}(P, Q) \in D[X]$ . The even exponent is useful in Chapter 2 and later when we deal with signs.

**Notation 1.16. [Truncation]** Let  $Q = b_q X^q + \cdots + b_0 \in D[X]$ . We define for  $0 \leq i \leq q$ , the **truncation of  $Q$  at  $i$**  by

$$\text{Tru}_i(Q) = b_i X^i + \cdots + b_0.$$

The **set of truncations** of a non-zero polynomial  $Q \in D[Y_1, \dots, Y_k][X]$ , where  $Y_1, \dots, Y_k$  are parameters and  $X$  is the main variable, is the finite subset of  $D[Y_1, \dots, Y_k][X]$  defined by

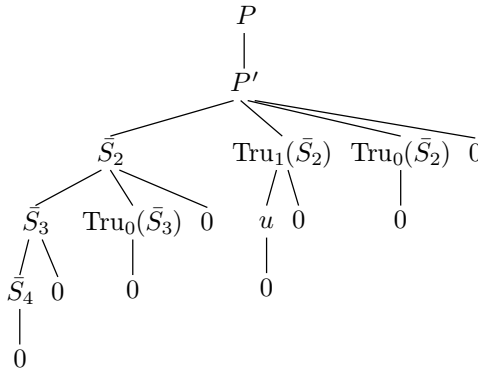
$$\text{Tru}(Q) = \begin{cases} \{Q\} & \text{if } \text{lcof}(Q) \in D \text{ or } \deg(Q) = 0, \\ \{Q\} \cup \text{Tru}(\text{Tru}_{\deg_X(Q)-1}(Q)) & \text{otherwise.} \end{cases}$$

The **tree of possible signed pseudo-remainder sequences** of two polynomials  $P, Q \in D[Y_1, \dots, Y_k][X]$ , denoted  $\text{TRems}(P, Q)$  is a tree whose root  $R$  contains  $P$ . The children of the root contain the elements of the set of truncations of  $Q$ . Each node  $N$  contains a polynomial  $\text{Pol}(N) \in D[Y_1, \dots, Y_k][X]$ . A node  $N$  is a leaf if  $\text{Pol}(N) = 0$ . If  $N$  is not a leaf, the children of  $N$  contain the truncations of  $-\text{PRem}(\text{Pol}(p(N)), \text{Pol}(N))$  where  $p(N)$  is the parent of  $N$ .  $\square$

*Example 1.17.* As in Example 1.15, we consider  $P = X^4 + aX^2 + bX + c$  and its derivative  $P' = 4X^3 + 2aX + b$ . Denoting

$$\begin{aligned} \bar{S}_2 &= -\text{PRem}(P, P') \\ &= -8aX^2 - 12bX - 16c, \\ \bar{S}_3 &= -\text{PRem}(P', \bar{S}_2) \\ &= 64((8ac - 9b^2 - 2a^3)X - b(12c + a^2)), \\ \bar{S}_4 &= -\text{PRem}(\bar{S}_3, \bar{S}_2) \\ &= 16384a^2(256c^3 - 128a^2c^2 + 144ab^2c + 16a^4c - 27b^4 - 4a^3b^2), \\ u &= -\text{PRem}(P', (\bar{S}_2)) \\ &= 768b(-27b^4 + 72acb^2 + 256c^3) \end{aligned}$$

the tree  $\text{TRems}(P, P')$  is the following.



Define

$$\begin{aligned} s &= 8ac - 9b^2 - 2a^3, \\ t &= -b(12c + a^2) \\ \delta &= 256c^3 - 128a^2c^2 + 144ab^2c + 16a^4c - 27b^4 - 4a^3b^2. \end{aligned}$$

The leftmost path in the tree going from the root to a leaf, namely the path  $P, P', S_2, S_3, S_4, 0$  can be understood as follows: if  $(a, b, c) \in C^3$  are such that the degree of the polynomials in the remainder sequence of  $P$  and  $P'$  are 4, 3, 2, 1, 0, i.e. when  $a \neq 0, s \neq 0, \delta \neq 0$  (getting rid of obviously irrelevant factors), then the signed remainder sequence of  $P = X^4 + aX^2 + bX + c$  and  $P'$  is proportional (up to non-zero squares of elements in  $C$ ) to  $P, P', \bar{S}_2, \bar{S}_3, \bar{S}_4, \square$

**Notation 1.18. [Degree]** For a specialization of  $Y = (Y_1, \dots, Y_k)$  to  $y \in C^k$ , and  $Q \in D[Y_1, \dots, Y_k][X]$ , we denote the polynomial in  $C[X]$  obtained by substituting  $y$  for  $Y$  by  $Q_y$ . Given  $\mathcal{Q} \subset D[Y_1, \dots, Y_k][X]$ , we define  $\mathcal{Q}_y \subset C[X]$  as  $\{Q_y \mid Q \in \mathcal{Q}\}$ .

Let  $Q = b_q X^q + \dots + b_0 \in D[Y_1, \dots, Y_k][X]$ . We define the basic formula  $\text{deg}_X(Q) = i$  as

$$\begin{cases} b_q = 0 \wedge \dots \wedge b_{i+1} = 0 \wedge b_i \neq 0 & \text{when } 0 \leq i < q, \\ b_q \neq 0 & \text{when } i = q, \\ b_q = 0 \wedge \dots \wedge b_0 = 0 & \text{when } i = -\infty, \end{cases}$$

so that the sets  $\text{Reali}(\text{deg}_X(Q) = i)$  partition  $C^k$  and  $y \in \text{Reali}(\text{deg}_X(Q) = i)$  if and only if  $\text{deg}(Q_y) = i$ .

Note that  $\text{PRem}(P_y, Q_y) = \text{PRem}(P, \text{Tru}_i(Q))_y$  where  $\text{deg}_X(Q_y) = i$ .

Given a leaf  $L$  of  $\text{TRems}(P, Q)$ , we denote by  $\mathcal{B}_L$  the unique path from the root of  $\text{TRems}(P, Q)$  to the leaf  $L$ . If  $N$  is a node in  $\mathcal{B}_L$  which is not a leaf, we denote by  $c(N)$  the unique child of  $N$  in  $\mathcal{B}_L$ . We denote by  $\mathcal{C}_L$  the basic formula

$$\bigwedge_{N \in \mathcal{B}_L, N \neq R} \text{deg}_X(-\text{PRem}(\text{Pol}(p(N)), \text{Pol}(N))) = \text{deg}_X(\text{Pol}(c(N))) \wedge \text{deg}_X(Q) = \text{deg}_X(\text{Pol}(c(R))) \wedge$$

□

It is clear from the definitions, since the remainder and pseudo-remainder of two polynomials in  $C[X]$  are equal up to a square, that

**Lemma 1.19.** *The  $\text{Reali}(\mathcal{C}_L)$  partition  $C^k$ . Moreover,  $y \in \text{Reali}(\mathcal{C}_L)$  implies that the signed remainder sequence of  $P_y$  and  $Q_y$  is proportional (up to a square) to the sequence of polynomials  $\text{Pol}(N)_y$  in the nodes along the path  $\mathcal{B}_L$  leading to  $L$ . In particular,  $\text{Pol}(p(L))_y$  is  $\text{gcd}(P_y, Q_y)$ .*

We will now define the set of possible greatest common divisors of a family  $\mathcal{P} \subset D[Y_1, \dots, Y_k][X]$ , called  $\text{posgcd}(\mathcal{P})$ , which is a finite set containing all the possible greatest common divisors of  $\mathcal{P}_y$  which can occur as  $y$  ranges over  $C^k$ . We define it as a set of pairs  $(G, \mathcal{C})$  where  $G \in D[Y_1, \dots, Y_k][X]$  and  $\mathcal{C}$  is a basic formula with coefficients in  $D$  so that for each pair  $(G, \mathcal{C})$ ,  $y \in \text{Reali}(\mathcal{C})$  implies  $\text{gcd}(\mathcal{P}_y) = G_y$ . More precisely, we shall make the definition so that the following lemma is true:

**Lemma 1.20.** *For all  $y \in C^k$ , there exists one and only one  $(G, \mathcal{C}) \in \text{posgcd}(\mathcal{P})$  such that  $y \in \text{Reali}(\mathcal{C})$ . Moreover,  $y \in \text{Reali}(\mathcal{C})$  implies that  $G_y$  is a greatest common divisor of  $\mathcal{P}_y$ .*

The set of possible greatest common divisors of a finite family of elements of  $\mathbf{K}[Y_1, \dots, Y_k][X]$  is defined recursively on the number of elements of the family by

$$\begin{aligned} \text{posgcd}(\emptyset) &= \{(0, 1 \neq 0)\} \\ \text{posgcd}(\mathcal{P} \cup \{P\}) &= \{(\text{Pol}(p(L)), \mathcal{C} \wedge \mathcal{C}_L) \mid (Q, \mathcal{C}) \in \text{posgcd}(\mathcal{P}) \\ &\quad \text{and } L \text{ is a leaf of } \text{TRems}(P, Q)\}. \end{aligned}$$

It is clear from the definitions and Lemma 1.19 that Lemma 1.20 holds.

*Example 1.21.* Returning to Example 1.17, and using the corresponding notation, the elements of  $\text{posgcd}(P, P')$  are (after removing obviously irrelevant factors),

$$\begin{aligned} (\bar{S}_4, a \neq 0 \wedge s \neq 0 \wedge \delta \neq 0), \\ (\bar{S}_3, a \neq 0 \wedge s \neq 0 \wedge \delta = 0), \\ (\text{Tru}_0(\bar{S}_3), a \neq 0 \wedge s = 0 \wedge t \neq 0), \\ (\bar{S}_2, a \neq 0 \wedge s = t = 0), \\ (u, a = 0 \wedge b \neq 0 \wedge u \neq 0), \\ (\text{Tru}_1(\bar{S}_2), a = 0 \wedge b \neq 0 \wedge u = 0), \\ (\text{Tru}_0(\bar{S}_2), a = b = 0 \wedge c \neq 0), \\ (P', a = b = c = 0). \end{aligned}$$

The first pair, which corresponds to the leftmost leaf of  $\text{TRems}(P, P')$  can be read as: if  $a \neq 0$ ,  $s \neq 0$ , and  $\delta \neq 0$  (i.e. if the degrees of the polynomials in the remainder sequence are 4, 3, 2, 1, 0), then  $\text{gcd}(P, P') = \bar{S}_4$ . The second pair, which corresponds to the next leaf (going left to right) means that if  $a \neq 0$ ,  $s \neq 0$ , and  $\delta = 0$  (i.e. if the degrees of the polynomials in the remainder sequence are 4, 3, 2, 1), then  $\text{gcd}(P, P') = \bar{S}_3$ .

If  $P = X^4 + aX^2 + bX + c$ , the projection of

$$\{(a, b, c, x) \in C^4 \mid P(x) = P'(x) = 0\}$$

to  $C^3$  is the set of polynomials (where a polynomial is identified with its coefficients  $(a, b, c)$ ) for which  $\deg(\text{gcd}(P, P')) \geq 1$ . Therefore, the formula  $\exists x P(x) = P'(x) = 0$  is equivalent to the formula

$$\begin{aligned} &(a \neq 0 \wedge s \neq 0 \wedge \delta = 0) \\ \vee &(a \neq 0 \wedge s = t = 0) \\ \vee &(a = 0 \wedge b \neq 0 \wedge u = 0) \\ \vee &(a = b = c = 0). \end{aligned}$$

□



The proof of the following projection theorem is based on the preceding constructions of possible gcd.

**Theorem 1.22. [Projection theorem for constructible sets]** *Given a constructible set in  $C^{k+1}$  defined by polynomials with coefficients in  $D$ , its projection to  $C^k$  is a constructible set defined by polynomials with coefficients in  $D$ .*

**Proof:** Since every constructible set is a finite union of basic constructible sets it is sufficient to prove that the projection of a basic constructible set is constructible. Suppose that the basic constructible set  $S$  in  $C^{k+1}$  is

$$\{(y, x) \in C^k \times C \mid \bigwedge_{P \in \mathcal{P}} P(y, x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(y, x) \neq 0\}$$

with  $\mathcal{P}$  and  $\mathcal{Q}$  finite subsets of  $D[Y_1, \dots, Y_k, X]$ .

Let

$$\mathcal{L} = \text{posgcd}(\{P \mid \exists \mathcal{C} \quad (P, \mathcal{C}) \in \text{posgcd}(\mathcal{P})\} \cup \{ \prod_{Q \in \mathcal{Q}} Q^d \})$$

where  $d$  is the least integer greater than the degree in  $X$  of any polynomial in  $\mathcal{P}$ .

For every  $(G, \mathcal{C}) \in \mathcal{L}$ , there exists a unique  $(G_1, \mathcal{C}_1) \in \text{posgcd}(\mathcal{P})$  with  $\mathcal{C}_1$  a conjunction of a subset of the atoms appearing in  $\mathcal{C}$ . Using Lemma 1.14, the projection of  $S$  on  $C^k$  is the union of the  $\text{Reali}(\mathcal{C} \wedge \deg_X(G) \neq \deg_X(G_1))$  for  $(G, \mathcal{C})$  in  $\mathcal{L}$ , and this is clearly a constructible set defined by polynomials with coefficients in  $D$ . □

**Exercise 1.9.**

- a) Find the conditions on  $(a, b, c)$  for  $P = aX^2 + bX + c$  and  $P' = 2aX + b$  to have a common root.
- b) Find the conditions on  $(a, b, c)$  for  $P = aX^2 + bX + c$  to have a root which is not a root of  $P'$ .

## 1.4 Quantifier Elimination and the Transfer Principle

Returning to logical terminology, Theorem 1.22 implies that the theory of algebraically closed fields admits quantifier elimination in the language of fields, which is the following theorem.

**Theorem 1.23. [Quantifier Elimination over Algebraically Closed Fields]** *Let  $\Phi(Y_1, \dots, Y_\ell)$  be a formula in the language of fields with free variables  $\{Y_1, \dots, Y_\ell\}$ , and coefficients in a subring  $D$  of the algebraically closed field  $C$ . Then there is a quantifier free formula  $\Psi(Y_1, \dots, Y_\ell)$  with coefficients in  $D$  which is  $C$ -equivalent to  $\Phi(Y_1, \dots, Y_\ell)$ .*

Notice that an example of quantifier elimination appears in Example 1.2.

The proof of the theorem is by induction on the number of quantifiers, using as base case the elimination of an existential quantifier which is given by Theorem 1.22.

**Proof of Theorem 1.23:** Given a formula  $\Theta(Y) = (\exists X) \mathcal{B}(X, Y)$ , where  $\mathcal{B}$  is a quantifier free formula whose atoms are equations and inequations involving polynomials in  $D[X, Y_1, \dots, Y_k]$ , Theorem 1.22 shows that there is a quantifier free formula  $\Xi(Y)$  with coefficients in  $D$  that is equivalent to  $\Theta(Y)$ , since  $\text{Reali}(\Theta(Y), C^k)$ , which is the projection of the constructible set  $\text{Reali}(\mathcal{B}(X, Y), C^{k+1})$ , is constructible, and constructible sets are realizations of quantifier free formulas. Since  $(\forall X) \Phi$  is equivalent to  $\neg((\exists X) \neg(\Phi))$ , the theorem immediately follows by induction on the number of quantifiers.  $\square$

**Corollary 1.24.** *Let  $\Phi(Y)$  be a formula in the language of fields with coefficients in  $C$ . The set  $\{y \in C^k \mid \Phi(y)\}$  is constructible.*

**Corollary 1.25.** *A subset of  $C$  defined by a formula in the language of fields with coefficients in  $C$  is a finite set or the complement of a finite set.*

**Proof:** By Corollary 1.24, a subset of  $C$  defined by a formula in the language of fields with coefficients in  $C$  is constructible, and this is a finite set or the complement of a finite set by Exercise 1.2.  $\square$

**Exercise 1.10.** Prove that the sets  $\mathbb{N}$  and  $\mathbb{Z}$  are not constructible subsets of  $\mathbb{C}$ . Prove that the sets  $\mathbb{N}$  and  $\mathbb{Z}$  cannot be defined inside  $\mathbb{C}$  by a formula of the language of fields with coefficients in  $\mathbb{C}$ .

Theorem 1.23 easily implies the following theorem, known as the transfer principle for algebraically closed fields. It is also called the Lefschetz Principle.

**Theorem 1.26. [Lefschetz principle]** *Suppose that  $C'$  is an algebraically closed field which contains the algebraically closed field  $C$ . If  $\Phi$  is a sentence in the language of fields with coefficients in  $C$ , then it is true in  $C$  if and only if it is true in  $C'$ .*

**Proof:** By Theorem 1.23, there is a quantifier free formula  $\Psi$  which is  $C$ -equivalent to  $\Phi$ . It follows from the proof of Theorem 1.22 that  $\Psi$  is  $C'$ -equivalent to  $\Phi$  as well. Notice, too, that since  $\Psi$  is a sentence,  $\Psi$  is a boolean combination of atoms of the form  $c = 0$  or  $c \neq 0$ , where  $c \in C$ . Clearly,  $\Psi$  is true in  $C$  if and only if it is true in  $C'$ .  $\square$

The **characteristic of a field  $K$**  is a prime number  $p$  if  $K$  contains  $\mathbb{Z}/p\mathbb{Z}$  and 0 if  $K$  contains  $\mathbb{Q}$ . The meaning of Lefschetz principle is essentially that a sentence is true in an algebraic closed field if and only if it is true in any other algebraic closed field of the same characteristic.

Let  $C$  denote an algebraically closed field and  $C'$  an algebraically closed field containing  $C$ .

Given a constructible set  $S$  in  $C^k$ , the **extension of  $S$  to  $C'$** , denoted  $\text{Ext}(S, C')$  is the constructible subset of  $C'^k$  defined by a quantifier free formula that defines  $S$ .

The following proposition is an easy consequence of Theorem 1.26.

**Proposition 1.27.** *Given a constructible set  $S$  in  $C^k$ , the set  $\text{Ext}(S, C')$  is well defined (i.e. it only depends on the set  $S$  and not on the quantifier free formula chosen to describe it).*

*The operation  $S \rightarrow \text{Ext}(S, C')$  preserves the boolean operations (finite intersection, finite union and complementation).*

*If  $S \subset T$ , then  $\text{Ext}(S, C') \subset \text{Ext}(T, C')$ , where  $T$  is a constructible set in  $C^k$ .*

**Exercise 1.11.** Prove proposition 1.27.

**Exercise 1.12.** Show that if  $S$  is a finite constructible subset of  $C^k$ , then  $\text{Ext}(S, C')$  is equal to  $S$ . (Hint: write a formula describing  $S$ ).

## 1.5 Bibliographical Notes

Lefschetz's principle (Theorem 1.26) is stated without proof in [105]. Indications for a proof of quantifier elimination over algebraically closed fields (Theorem 1.23) are given in [156] (Remark 16).

---

## Real Closed Fields

Real closed fields are fields which share the algebraic properties of the field of real numbers. In Section 2.1, we define ordered, real and real closed fields and state some of their basic properties. Section 2.2 is devoted to real root counting. In Section 2.3 we define semi-algebraic sets and prove that the projection of an algebraic set is semi-algebraic. The main technique used is a parametric version of real root counting algorithm described in the second section. In Section 2.4, we prove that the projection of a semi-algebraic set is semi-algebraic, by a similar method. Section 2.5 is devoted to several applications of the projection theorem, of logical and geometric nature. In Section 2.6, an important example of a non-archimedean real closed field is described: the field of Puiseux series.

### 2.1 Ordered, Real and Real Closed Fields

Before defining ordered fields, we prove a few useful properties of fields of characteristic zero.

Let  $K$  be a field of characteristic zero. The **derivative** of a polynomial

$$P = a_p X^p + \dots + a_i X^i + \dots + a_0 \in K[X]$$

is denoted  $P'$  with

$$P' = p a_p X^{p-1} + \dots + i a_i X^{i-1} + \dots + a_1.$$

The  $i$ -th derivative of  $P$ ,  $P^{(i)}$ , is defined inductively by  $P^{(i)} = \left(P^{(i-1)}\right)'$ . It is immediate to verify that

$$\begin{aligned} (P + Q)' &= P' + Q', \\ (PQ)' &= P'Q + PQ'. \end{aligned}$$

Taylor's formula holds:

**Proposition 2.1. [Taylor's formula]** *Let  $K$  be a field of characteristic zero,*

$$P = a_p X^p + \dots + a_i X^i + \dots + a_0 \in K[X] \text{ and } x \in K.$$

Then,

$$P = \sum_{i=0}^{\deg(P)} \frac{P^{(i)}(x)}{i!} (X-x)^i.$$

**Proof:** We prove Taylor's formula holds for monomials  $X^p$  by induction on  $p$ . The claim is clearly true if  $p=0$ . Suppose that Taylor's formula holds for  $p-1$ :

$$X^{p-1} = \sum_{i=0}^{p-1} \frac{(p-1)!}{(p-1-i)! i!} x^{p-1-i} (X-x)^i.$$

Then, since  $X = x + (X-x)$ ,

$$\begin{aligned} X^p &= (x + (X-x)) \sum_{i=0}^{p-1} \frac{(p-1)!}{(p-1-i)! i!} x^{p-1-i} (X-x)^i \\ &= \sum_{i=0}^p \frac{p!}{(p-i)! i!} x^{p-i} (X-x)^i \end{aligned}$$

since

$$\frac{p!}{(p-i)! i!} = \frac{(p-1)!}{(p-i)! (i-1)!} + \frac{p!}{(p-1-i)! (i-1)!}.$$

Hence, Taylor's formula is valid for any polynomial using the linearity of derivation.  $\square$

Let  $x \in K$  and  $P \in K[X]$ . The **multiplicity** of  $x$  as a root of  $P$  is the natural number  $\mu$  such that there exists  $Q \in K[X]$  with  $P = (X-x)^\mu Q(X)$  and  $Q(x) \neq 0$ . Note that if  $x$  is not a root of  $P$ , the multiplicity of  $x$  as a root of  $P$  is equal to 0.

**Lemma 2.2.** *Let  $K$  be a field of characteristic zero. The element  $x \in K$  is a root of  $P \in K[X]$  of multiplicity  $\mu$  if and only if*

$$P^{(\mu)}(x) \neq 0, P^{(\mu-1)}(x) = \dots = P(x) = P'(x) = 0.$$

**Proof:** Suppose that  $P = (X-x)^\mu Q$  and  $Q(x) \neq 0$ . It is clear that  $P(x) = 0$ . The proof of the claim is by induction on the degree of  $P$ . The claim is obviously true for  $\deg(P) = 1$ . Suppose that the claim is true for every polynomial of degree  $< d$ . Since

$$P' = (X-x)^{\mu-1} (\mu Q + (X-x) Q'),$$

and  $\mu Q(x) \neq 0$ , by induction hypothesis,

$$P'(x) = \dots = P^{(\mu-1)}(x) = 0, P^{(\mu)}(x) \neq 0.$$

Conversely suppose that

$$P(x) = P'(x) = \dots = P^{(\mu-1)}(x) = 0, P^{(\mu)}(x) \neq 0.$$

By Proposition 2.1 (Taylor’s formula) at  $x$ ,  $P = (X - x)^\mu Q$ , with

$$Q(x) = P^{(\mu)}(x)/\mu! \neq 0. \quad \square$$

A polynomial  $P \in K[X]$  is **separable** if the greatest common divisor of  $P$  and  $P'$  is an element of  $K \setminus \{0\}$ . A polynomial  $P$  is **square-free** if there is no non-constant polynomial  $A \in K[X]$  such that  $A^2$  divides  $P$ .

**Exercise 2.1.** Prove that  $P \in K[X]$  is separable if and only if  $P$  has no multiple root in  $C$ , where  $C$  is an algebraically closed field containing  $K$ . If the characteristic of  $K$  is 0, prove that  $P \in K[X]$  is separable if and only if  $P$  is square-free.

A **partially ordered set**  $(A, \preceq)$  is a set  $A$ , together with a binary relation  $\preceq$  that satisfies:

- $\preceq$  is transitive, i.e.  $a \preceq b$  and  $a \preceq c \Rightarrow a \preceq c$ ,
- $\preceq$  is reflexive, i.e.  $a \preceq a$ ,
- $\preceq$  is anti-symmetric, i.e.  $a \preceq b$  and  $b \preceq a \Rightarrow a = b$ .

A standard example of a partially ordered set is the power set

$$2^A = \{B \mid B \subseteq A\},$$

the binary relation being the inclusion between subsets of  $A$ .

A **totally ordered set** is a partially ordered set  $(A, \leq)$  with the additional property that every two elements  $a, b \in A$  are comparable, i.e.  $a \leq b$  or  $b \leq a$ . In a totally ordered set,  $a < b$  stands for  $a \leq b$ ,  $a \neq b$ , and  $a \geq b$  (resp.  $a > b$ ) for  $b \leq a$  (resp.  $b < a$ ).

An **ordered ring**  $(A, \leq)$  is a ring,  $A$ , together with a total order,  $\leq$ , that satisfies:

$$\begin{aligned} x \leq y &\Rightarrow x + z \leq y + z \\ 0 \leq x, 0 \leq y &\Rightarrow 0 \leq xy. \end{aligned}$$

An **ordered field**  $(F, \leq)$  is a field,  $F$ , which is an ordered ring.

An ordered ring  $(A, \leq)$  is contained in an ordered field  $(F, \leq)$  if  $A \subset F$  and the inclusion is order preserving. Note that the ordered ring  $(A, \leq)$  is necessarily an ordered integral domain.

**Exercise 2.2.** Prove that in an ordered field  $-1 < 0$ .

Prove that an ordered field has characteristic zero.

Prove the law of trichotomy in an ordered field: for every  $a$  in the field, exactly one of  $a < 0$ ,  $a = 0$ ,  $a > 0$  holds.

**Notation 2.3. [Sign]** The **sign** of an element  $a$  in ordered field  $(F, \leq)$  is defined by

$$\begin{cases} \text{sign}(a) = 0 & \text{if } a = 0, \\ \text{sign}(a) = 1 & \text{if } a > 0, \\ \text{sign}(a) = -1 & \text{if } a < 0. \end{cases}$$

When  $a > 0$  we say  $a$  is positive, and when  $a < 0$  we say  $a$  is negative.

The **absolute value**  $|a|$  of  $a$  is the maximum of  $a$  and  $-a$  and is non-negative.  $\square$

The fields  $\mathbb{Q}$  and  $\mathbb{R}$  with their natural order are familiar examples of ordered fields.

**Exercise 2.3.** Show that it is not possible to order the field of complex numbers  $\mathbb{C}$  so that it becomes an ordered field.

In an ordered field, the value at  $x$  of a polynomial has the sign of its leading monomial for  $x$  sufficiently large. More precisely,

**Proposition 2.4.** *Let  $P = a_p X^p + \dots + a_0$ ,  $a_p \neq 0$ , be a polynomial with coefficients in an ordered field  $F$ . If  $|x|$  is bigger than  $2 \sum_{0 \leq i \leq p} \frac{|a_i|}{|a_p|}$ , then  $P(x)$  and  $a_p x^p$  have the same sign.*

**Proof:** Suppose that

$$|x| > 2 \sum_{0 \leq i \leq p} \left| \frac{a_i}{a_p} \right|,$$

which implies  $|x| > 2$ . Since

$$\begin{aligned} \frac{P(x)}{a_p x^p} &= 1 + \sum_{0 \leq i \leq p-1} \frac{a_i}{a_p} x^{i-p}, \\ \frac{P(x)}{a_p x^p} &\geq 1 - \left( \sum_{0 \leq i \leq p-1} \frac{|a_i|}{|a_p|} |x|^{i-p} \right) \\ &\geq 1 - \left( \sum_{0 \leq i \leq p} \frac{|a_i|}{|a_p|} \right) (|x|^{-1} + |x|^{-2} + \dots + |x|^{-p}) \\ &\geq 1 - \frac{1}{2} (1 + |x|^{-1} + \dots + |x|^{-p+1}) \\ &= 1 - \frac{1}{2} \left( \frac{1 - |x|^{-p}}{1 - |x|^{-1}} \right) > 0. \end{aligned}$$

$\square$

We now examine a particular way to order the field of rational functions  $\mathbb{R}(X)$ .

For this purpose, we need a definition: Let  $F \subset F'$  be two ordered fields. The element  $f \in F'$  is **infinitesimal over  $F$**  if it is a positive element smaller than any positive  $f \in F$ . The element  $f \in F'$  is **unbounded over  $F$**  if it is a positive element greater than any positive  $f \in F$ .

**Notation 2.5. [Order  $0_+$ ]** Let  $F$  be an ordered field and  $\varepsilon$  a variable. There is one and only one order on  $F(\varepsilon)$ , denoted  $0_+$ , such that  $\varepsilon$  is infinitesimal over  $F$ . If

$$P(\varepsilon) = a_p \varepsilon^p + a_{p-1} \varepsilon^{p-1} + \dots + a_{m+1} \varepsilon^{m+1} + a_m \varepsilon^m$$

with  $a_m \neq 0$ , then  $P(\varepsilon) > 0$  in this order if and only if  $a_m > 0$ . If  $P(\varepsilon)/Q(\varepsilon) \in F(\varepsilon)$ ,  $P(\varepsilon)/Q(\varepsilon) > 0$  if and only if  $P(\varepsilon)Q(\varepsilon) > 0$ .

Note that the field  $F(\varepsilon)$  with this order contains infinitesimal elements over  $F$ , such as  $\varepsilon$ . The field also contains elements which are unbounded over  $F$  such as  $1/\varepsilon$ . □

**Exercise 2.4.** Show that  $0_+$  is an order on  $F(\varepsilon)$  and that it is the only order in which  $\varepsilon$  is infinitesimal over  $F$ .

We define now a cone of a field, which should be thought of as a set of non-negative elements. A **cone** of the field  $F$  is a subset  $\mathcal{C}$  of  $F$  such that:

$$\begin{aligned} x \in \mathcal{C}, y \in \mathcal{C} &\Rightarrow x + y \in \mathcal{C} \\ x \in \mathcal{C}, y \in \mathcal{C} &\Rightarrow xy \in \mathcal{C} \\ x \in F &\Rightarrow x^2 \in \mathcal{C}. \end{aligned}$$

The cone  $\mathcal{C}$  is **proper** if in addition  $-1 \notin \mathcal{C}$ .

Let  $(F, \leq)$  be an ordered field. The subset  $\mathcal{C} = \{x \in F \mid x \geq 0\}$  is a cone, the **positive cone** of  $(F, \leq)$ .

**Proposition 2.6.** *Let  $(F, \leq)$  be an ordered field. The positive cone  $\mathcal{C}$  of  $(F, \leq)$  is a proper cone that satisfies  $\mathcal{C} \cup -\mathcal{C} = F$ . Conversely, if  $\mathcal{C}$  is a proper cone of a field  $F$  that satisfies  $\mathcal{C} \cup -\mathcal{C} = F$ , then  $F$  is ordered by  $x \leq y \Leftrightarrow y - x \in \mathcal{C}$ .*

**Exercise 2.5.** Prove Proposition 2.6.

Let  $K$  be a field. We denote by  $K^{(2)}$  the set of squares of elements of  $K$  and by  $\sum K^{(2)}$  the set of **sums of squares of elements of  $K$** . Clearly,  $\sum K^{(2)}$  is a cone contained in every cone of  $K$ .

A field  $K$  is a **real field** if  $-1 \notin \sum K^{(2)}$ .

**Exercise 2.6.** Prove that a real field has characteristic 0.

Show that the field  $\mathbb{C}$  of complex numbers is not a real field.

Show that an ordered field is a real field.

Real fields can be characterized as follows.

**Theorem 2.7.** *Let  $F$  be a field. Then the following properties are equivalent*

- a)  $F$  is real.
- b)  $F$  has a proper cone.
- c)  $F$  can be ordered.
- d) For every  $x_1, \dots, x_n$  in  $F$ ,  $\sum_{i=1}^n x_i^2 = 0 \Rightarrow x_1 = \dots = x_n = 0$ .

The proof of Theorem 2.7 uses the following proposition.

**Proposition 2.8.** *Let  $\mathcal{C}$  be a proper cone of  $F$ ,  $\mathcal{C}$  is contained in the positive cone for some order on  $F$ .*



The proof of Proposition 2.8 relies on the following lemma.

**Lemma 2.9.** *Let  $\mathcal{C}$  be a proper cone of  $F$ . If  $-a \notin \mathcal{C}$ , then*

$$\mathcal{C}[a] = \{x + ay \mid x, y \in \mathcal{C}\}$$

*is a proper cone of  $F$ .*

**Proof:** Suppose  $-1 = x + ay$  with  $x, y \in \mathcal{C}$ . If  $y = 0$  we have  $-1 \in \mathcal{C}$  which is impossible. If  $y \neq 0$  then  $-a = (1/y^2)y(1+x) \in \mathcal{C}$ , which is also impossible.  $\square$

**Proof of Proposition 2.8:** Since the union of a chain of proper cones is a proper cone, Zorn's lemma implies the existence of a maximal proper cone  $\bar{\mathcal{C}}$  which contains  $\mathcal{C}$ . It is then sufficient to show that  $\bar{\mathcal{C}} \cup -\bar{\mathcal{C}} = F$ , and to define  $x \leq y$  by  $y - x \in \bar{\mathcal{C}}$ . Suppose that  $-a \notin \bar{\mathcal{C}}$ . By Lemma 2.9,  $\bar{\mathcal{C}}[a]$  is a proper cone and thus, by the maximality of  $\bar{\mathcal{C}}$ ,  $\bar{\mathcal{C}} = \bar{\mathcal{C}}[a]$  and thus  $a \in \bar{\mathcal{C}}$ .  $\square$

**Proof of Theorem 2.7:**

- $a) \Rightarrow b)$  since in a real field  $F$ ,  $\sum F^{(2)}$  is a proper cone.
- $b) \Rightarrow c)$  by Proposition 2.8.
- $c) \Rightarrow d)$  since in an ordered field, if  $x_1 \neq 0$  then  $\sum_{i=1}^n x_i^2 \geq x_1^2 > 0$ .
- $d) \Rightarrow a)$ , since in a field  $0 \neq 1$ , so 4 implies that  $1 + \sum_{i=1}^n x_i^2 = 0$  is impossible.  $\square$

A **real closed field**  $R$  is an ordered field whose positive cone is the set of squares  $R^{(2)}$  and such that every polynomial in  $R[X]$  of odd degree has a root in  $R$ .

Note that the condition that the positive cone of a real closed field  $R$  is  $R^{(2)}$  means that  $R$  has a unique order as an ordered field, since the positive cone of an order contains necessarily  $R^{(2)}$ .

*Example 2.10.* The field  $\mathbb{R}$  of real numbers is of course real closed. The **real algebraic numbers**, i.e. those real numbers that satisfy an equation with integer coefficients, form a real closed field denoted  $\mathbb{R}_{\text{alg}}$  (see Exercise 2.11)  $\square$

A field  $R$  has the **intermediate value property** if  $R$  is an ordered field such that, for any  $P \in R[X]$ , if there exist  $a \in R, b \in R, a < b$  such that  $P(a)P(b) < 0$ , there exists  $x \in (a, b)$  such that  $P(x) = 0$ .

Real closed fields are characterized as follows.

**Theorem 2.11.** *If  $R$  is a field then the following properties are equivalent:*

- a)  $R$  is real closed.*
- b)  $R[i] = R[T]/(T^2 + 1)$  is an algebraically closed field.*
- c)  $R$  has the intermediate value property.*
- d)  $R$  is a real field that has no non-trivial real algebraic extension, that is there is no real field  $R_1$  that is algebraic over  $R$  and different from  $R$ .*

The following classical definitions and results about symmetric polynomials are used in the proof of Theorem 2.11.

Let  $K$  be a field. A polynomial  $Q(X_1, \dots, X_k) \in K[X_1, \dots, X_k]$  is **symmetric** if for every permutation  $\sigma$  of  $\{1, \dots, k\}$ ,

$$Q(X_{\sigma(1)}, \dots, X_{\sigma(k)}) = Q(X_1, \dots, X_k).$$

**Exercise 2.7.** Denote by  $\mathcal{S}_k$  the group of permutations of  $\{1, \dots, k\}$ .

If  $X^\alpha = X_1^{\alpha_1} \dots X_k^{\alpha_k}$ , denote  $X_\sigma^\alpha = X_{\sigma(1)}^{\alpha_1} \dots X_{\sigma(k)}^{\alpha_k}$  and  $M_\alpha = \sum_{\sigma \in \mathcal{S}_p} X_\sigma^\alpha$ . Prove that every symmetric polynomial can be written as a finite sum  $\sum c_\alpha M_\alpha$ .

For  $i = 1, \dots, k$ , the  **$i$ -th elementary symmetric function** is

$$E_i = \sum_{1 \leq j_1 < \dots < j_i \leq k} X_{j_1} \dots X_{j_i}.$$

Elementary symmetric functions are related to coefficients of polynomials as follows.

**Lemma 2.12.** *Let  $X_1, \dots, X_k$  be elements of a field  $K$  and*

$$P = (X - X_1) \dots (X - X_k) = X^k + C_1 X^{k-1} + \dots + C_k,$$

*then  $C_i = (-1)^i E_i$ .*

**Proof:** Identify the coefficient of  $X^i$  on both sides of

$$(X - X_1) \dots (X - X_k) = X^k + C_1 X^{k-1} + \dots + C_k. \quad \square$$

**Proposition 2.13.** *Let  $K$  be a field and let*

$$Q(X_1, \dots, X_k) \in K[X_1, \dots, X_k]$$

*be symmetric. There exists a polynomial*

$$R(T_1, \dots, T_k) \in K[T_1, \dots, T_k]$$

*such that  $Q(X_1, \dots, X_k) = R(E_1, \dots, E_k)$ .*

The proof of Proposition 2.13 uses the notion of graded lexicographical ordering. We define first the lexicographical ordering, which is the order of the dictionary and will be used at several places in the book.

We denote by  $\mathcal{M}_k$  the set of monomials in  $k$  variables. Note that  $\mathcal{M}_k$  can be identified with  $\mathbb{N}^k$  defining  $X^\alpha = X_1^{\alpha_1} \dots X_k^{\alpha_k}$ .

**Definition 2.14. [Lexicographical ordering]** Let  $(B, <)$  be a totally ordered set. The **lexicographical ordering**,  $<_{\text{lex}}$ , on finite sequences of  $k$  elements of  $B$  is the total order  $<_{\text{lex}}$  defined by induction on  $k$  by

$$b <_{\text{lex}} b' \Leftrightarrow b < b'$$

$$(b_1, \dots, b_k) <_{\text{lex}} (b'_1, \dots, b'_k) \Leftrightarrow (b_1 < b'_1) \vee (b_1 = b'_1 \wedge (b_2, \dots, b_k) <_{\text{lex}} (b'_2, \dots, b'_k)).$$

We denote by  $\mathcal{M}_k$  the set of monomials in  $k$  variables  $X_1, \dots, X_k$ . Note that  $\mathcal{M}_k$  can be identified with  $\mathbb{N}^k$  defining  $X^\alpha = X_1^{\alpha_1} \cdots X_k^{\alpha_k}$ . Using this identification defines the lexicographical ordering  $<_{\text{lex}}$  on  $\mathcal{M}_k$ . In the lexicographical ordering,  $X_1 >_{\text{grlex}} \cdots >_{\text{grlex}} X_k$ . The smallest monomial with respect to the lexicographical ordering is 1, and the lexicographical ordering is compatible with multiplication. Note that the set of monomials less than or equal to a monomial  $X^\alpha$  in the lexicographical ordering maybe infinite.  $\square$

**Exercise 2.8.** Prove that a strictly decreasing sequence for the lexicographical ordering is necessarily finite. Hint: by induction on  $k$ .

**Definition 2.15. [Graded lexicographical ordering]** The **graded lexicographical ordering**,  $<_{\text{grlex}}$ , on the set of monomials in  $k$  variables  $\mathcal{M}_k$  is the total order  $X^\alpha <_{\text{grlex}} X^\beta$  defined by

$$X^\alpha <_{\text{grlex}} X^\beta \Leftrightarrow (\deg(X^\alpha) < \deg(X^\beta)) \vee (\deg(X^\alpha) = \deg(X^\beta) \wedge \alpha <_{\text{lex}} \beta)$$

with  $\alpha = (\alpha_1, \dots, \alpha_k), \beta = (\beta_1, \dots, \beta_k), X^\alpha = X_1^{\alpha_1} \cdots X_k^{\alpha_k}, X^\beta = X_1^{\beta_1} \cdots X_k^{\beta_k}$ .

In the graded lexicographical ordering above,  $X_1 >_{\text{grlex}} \cdots >_{\text{grlex}} X_k$ . The smallest monomial with respect to the graded lexicographical ordering is 1, and the graded lexicographical ordering is compatible with multiplication. Note that the set of monomials less than or equal to a monomial  $X^\alpha$  in the graded lexicographical ordering is finite.  $\square$

**Proof of Proposition 2.13:** Since  $Q(X_1, \dots, X_k)$  is symmetric, its leading monomial in the graded lexicographical ordering  $c_\alpha X^\alpha = c_\alpha X_1^{\alpha_1} \cdots X_k^{\alpha_k}$  satisfies  $\alpha_1 \geq \dots \geq \alpha_k$ . The leading monomial of  $c_\alpha E_1^{\alpha_1 - \alpha_2} \cdots E_{k-1}^{\alpha_{k-1} - \alpha_k} E_k^{\alpha_k}$  in the graded lexicographical ordering is also  $c_\alpha X^\alpha = c_\alpha X_1^{\alpha_1} \cdots X_k^{\alpha_k}$ .

Let  $Q_1 = Q(X_1, \dots, X_k) - c_\alpha E_1^{\alpha_1 - \alpha_2} \cdots E_{k-1}^{\alpha_{k-1} - \alpha_k} E_k^{\alpha_k}$ . If  $Q_1 = 0$ , the proof is over. Otherwise, the leading monomial with respect to the graded lexicographical ordering of  $Q_1$  is strictly smaller than  $X_1^{\alpha_1} \cdots X_k^{\alpha_k}$ , and it is possible to iterate the construction with  $Q_1$ . Since there is no infinite decreasing sequence of monomials for the graded lexicographical ordering, the claim follows.  $\square$

**Proposition 2.16.** Let  $P \in K[X]$ , of degree  $k$ , and  $x_1, \dots, x_k$  be the roots of  $P$  (counted with multiplicities) in an algebraically closed field  $\mathbb{C}$  containing  $K$ . If a polynomial  $Q(X_1, \dots, X_k) \in K[X_1, \dots, X_k]$  is symmetric, then  $Q(x_1, \dots, x_k) \in K$ .

**Proof:** Let  $e_i$ , for  $1 \leq i \leq k$ , denote the  $i$ -th elementary symmetric function evaluated at  $x_1, \dots, x_k$ . Since  $P \in K[X]$ , Lemma 2.12 gives  $e_i \in K$ . By Proposition 2.13, there exists  $R(T_1, \dots, T_k) \in K[T_1, \dots, T_k]$  such that

$$Q(X_1, \dots, X_k) = R(E_1, \dots, E_k).$$

Thus,  $Q(x_1, \dots, x_k) = R(e_1, \dots, e_k) \in K$ .  $\square$

With these preliminaries results, it is possible to prove Theorem 2.11.

**Proof of Theorem 2.11:**  $a) \Rightarrow b)$  Let  $P \in \mathbb{R}[X]$  a monic separable polynomial of degree  $p = 2^m n$  with  $n$  odd. We show by induction on  $m$  that  $P$  has a root in  $\mathbb{R}[i]$ .

If  $m = 0$ , then  $p$  is odd and  $P$  has a root in  $\mathbb{R}$ , since  $\mathbb{R}$  is real closed.

Denote by  $x_1, \dots, x_p$  the roots of  $P$  in an algebraically closed field  $\mathbb{C}$ . Let  $Z$  be a new indeterminate and  $Q(Z, Y)$  the monic polynomial having as roots the  $x_i + x_j + Z x_i x_j$  where  $i < j$ .

$$Q(Z, Y) = \prod_{i < j} (Y - (x_i + x_j + Z x_i x_j)).$$

The coefficients of  $Q(Z, Y)$  can be explicitly computed as polynomials of the coefficients of  $P$ , using Proposition 2.16, thus  $Q(Z, Y) \in \mathbb{R}[Z, Y]$ . The degree of  $Q(Z, Y)$  in  $Y$  and  $Z$  is  $p(p-1)/2$ .

Ordering lexicographically the couples  $(i, j)$ ,  $i < j$ , we define the discriminant of  $Q$  as

$$\begin{aligned} D(Z) &= \prod_{\substack{i < j, k < \ell \\ (i, j) < (k, \ell)}} ((x_i + x_j + Z x_i x_j) - (x_k + x_\ell + Z x_k x_\ell))^2 \\ &= \prod_{\substack{i < j, k < \ell \\ (i, j) < (k, \ell)}} (\alpha_{i, j, k, \ell} + Z \beta_{i, j, k, \ell})^2 \end{aligned}$$

where  $\alpha_{i, j, k, \ell} = (x_i + x_j - x_k + x_\ell)$ ,  $\beta_{i, j, k, \ell} = x_i x_j - x_k x_\ell$ . Note that by Proposition 2.16,  $D(Z) \in \mathbb{R}[Z]$ .

Since all the roots of  $P$  are distinct, we get the following implication

$$i < j, k < \ell, (i, j) < (k, \ell), x_i x_j = x_k x_\ell \Rightarrow x_i + x_j \neq x_k + x_\ell.$$

So every factor  $\alpha_{i, j, k, \ell} + Z \beta_{i, j, k, \ell}$  is nonzero. It follows that  $D(Z)$  is not identically zero.

Taking a value  $z \in \mathbb{N}$  such that  $D(z) \neq 0$ , the polynomial  $Q(z, Y)$  is a square free polynomial since all its roots are distinct.

We prove now that it is possible to express, for every  $1 \leq i < j \leq p$ ,  $x_i + x_j$  and  $x_i x_j$  rationally in terms of  $\gamma_{i, j} = x_i + x_j + z x_i x_j$ .

Indeed let

$$\begin{aligned} F(Z, Y) &= \partial Q / \partial Y(Z, Y) \\ &= \sum_{i < j} \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (Y - (x_k + x_\ell + Z x_k x_\ell)) \\ G(Z, Y) &= \sum_{i < j} (x_i + x_j) \left( \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (Y - (x_k + x_\ell + Z x_k x_\ell)) \right), \\ H(Z, Y) &= \sum_{i < j} x_i x_j \left( \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (Y - (x_k + x_\ell + Z x_k x_\ell)) \right). \end{aligned}$$

Note that by Proposition 2.16,  $f(Z, Y)$ ,  $G(Z, Y)$  and  $H(Z, Y)$  are elements of  $\mathbb{R}[Z, Y]$ .

Then, for every  $1 \leq i < j \leq p$ ,

$$\begin{aligned} F(z, \gamma_{i,j}) &= \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (\gamma_{i,j} - \gamma_{k,\ell}), \\ G(z, \gamma_{i,j}) &= (x_i + x_j) \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (\gamma_{i,j} - \gamma_{k,\ell}), \\ H(z, \gamma_{i,j}) &= (x_i x_j) \prod_{\substack{k < \ell \\ (k, \ell) \neq (i, j)}} (\gamma_{i,j} - \gamma_{k,\ell}). \end{aligned}$$

It follows that

$$\begin{aligned} x_i + x_j &= \frac{G(z, \gamma_{i,j})}{F(z, \gamma_{i,j})}, \\ x_i x_j &= \frac{H(z, \gamma_{i,j})}{F(z, \gamma_{i,j})}. \end{aligned}$$

In other words, the roots of the second degree polynomial

$$F(z, \gamma_{i,j})X^2 - G(z, \gamma_{i,j})X + H(z, \gamma_{i,j})$$

are roots of  $P$ .

The polynomial  $Q(z, Y)$  is of degree  $p(p-1)/2$ , i.e. of the form  $2^{m-1}n'$  with  $n'$  odd. By induction hypothesis, it has a root  $\gamma$  in  $\mathbf{R}[i]$ . Since the classical method for solving polynomials of degree 2 works in  $\mathbf{R}[i]$  when  $\mathbf{R}$  is real closed, the roots of the second degree polynomial

$$F(z, \gamma)X^2 - G(z, \gamma)X + H(z, \gamma)$$

are roots of  $P$  that belong to  $\mathbf{R}[i]$ . We have proved that the polynomial  $P$  has a root in  $\mathbf{R}[i]$ .

For  $P = a_p X^p + \dots + a_0 \in \mathbf{R}[i][X]$ , we write  $\bar{P} = \bar{a}_p X^p + \dots + \bar{a}_0$ . Since  $P\bar{P} \in \mathbf{R}[X]$ ,  $P\bar{P}$  has a root  $x$  in  $\mathbf{R}[i]$ . Thus  $P(x) = 0$  or  $\bar{P}(x) = 0$ . In the first case we are done and in the second,  $P(\bar{x}) = 0$ .

$b) \Rightarrow c)$  Since  $\mathbf{C} = \mathbf{R}[i]$  is algebraically closed,  $P$  factors into linear factors over  $\mathbf{C}$ . Since if  $c + id$  is a root of  $P$ ,  $c - id$  is also a root of  $P$ , the irreducible factors of  $P$  are linear or have the form

$$(X - c)^2 + d^2 = (X - c - id)(X - c + id), \quad d \neq 0.$$

If  $P(a)$  and  $P(b)$  have opposite signs, then  $Q(a)$  and  $Q(b)$  have opposite signs for some linear factor  $Q$  of  $P$ . Hence the root of  $Q$  is in  $(a, b)$ .

$c) \Rightarrow a)$  If  $y$  is positive,  $X^2 - y$  takes a negative value at 0 and a positive value for  $X$  big enough, by Proposition 2.4. Thus  $X^2 - y$  has a root, which is a square root of  $y$ . Similarly a polynomial of odd degree with coefficients in  $\mathbf{R}$  takes different signs for  $a$  positive and big enough and  $b$  negative and small enough, using Proposition 2.4 again. Thus it has a root in  $\mathbf{R}$ .

$b) \Rightarrow d)$  Since  $\mathbb{R}[i] = \mathbb{R}[T]/(T^2 + 1)$  is a field,  $T^2 + 1$  is irreducible over  $\mathbb{R}$ . Hence  $-1$  is not a square in  $\mathbb{R}$ . Moreover in  $\mathbb{R}$ , a sum of squares is still a square: let  $a, b \in \mathbb{R}$  and  $c, d \in \mathbb{R}$  such that  $a + ib = (c + id)^2$ ; then  $a^2 + b^2 = (c^2 + d^2)^2$ . This proves that  $\mathbb{R}$  is real. Finally, since the only irreducible polynomials of  $\mathbb{R}[X]$  of degree  $> 1$  are of the form

$$(X - c)^2 + d^2 = (X - c - id)(X - c + id), \quad d \neq 0,$$

and  $\mathbb{R}[X]/((X - c)^2 + d^2) = \mathbb{R}[i]$ , the only non-trivial algebraic extensions of  $\mathbb{R}$  is  $\mathbb{R}[i]$ , which is not real.

$d) \Rightarrow a)$  Suppose that  $a \in \mathbb{R}$ . If  $a$  is not a square in  $\mathbb{R}$ , then

$$\mathbb{R}[\sqrt{a}] = \mathbb{R}[X]/(X^2 - a)$$

is a non-trivial algebraic extension of  $\mathbb{R}$ , and thus  $\mathbb{R}[\sqrt{a}]$  is not real. Thus,

$$\begin{aligned} -1 &= \sum_{i=1}^n (x_i + \sqrt{a} y_i)^2 \\ -1 &= \sum_{i=1}^n x_i^2 + a \sum_{i=1}^n y_i^2 \in \mathbb{R}. \end{aligned}$$

Since  $\mathbb{R}$  is real,  $-1 \neq \sum_{i=1}^n x_i^2$  and thus  $y = \sum_{i=1}^n y_i^2 \neq 0$ . Hence,

$$\begin{aligned} -a &= \left( \sum_{i=1}^n y_i^2 \right)^{-1} \left( 1 + \sum_{i=1}^n x_i^2 \right) \\ &= \left( \sum_{i=1}^n \left( \frac{y_i}{y} \right)^2 \right) \left( 1 + \sum_{i=1}^n x_i^2 \right) \in \sum \mathbb{R}^{(2)}. \end{aligned}$$

This shows that  $\mathbb{R}^{(2)} \cup -\sum \mathbb{R}^{(2)} = \mathbb{R}$  and thus that there is only one possible order on  $\mathbb{R}$  with  $\mathbb{R}^{(2)} = \sum \mathbb{R}^{(2)}$  as positive cone.

It remains to show that if  $P \in \mathbb{R}[X]$  has odd degree then  $P$  has a root in  $\mathbb{R}$ . If this is not the case, let  $P$  be a polynomial of odd degree  $p > 1$  such that every polynomial of odd degree  $< p$  has a root in  $\mathbb{R}$ . Since a polynomial of odd degree has at least one odd irreducible factor, we assume without loss of generality that  $P$  is irreducible. The quotient  $\mathbb{R}[X]/(P)$  is a non-trivial algebraic extension of  $\mathbb{R}$  and hence  $-1 = \sum_{i=1}^n H_i^2 + PQ$  with  $\deg(H_i) < p$ . Since the term of highest degree in the expansion of  $\sum_{i=1}^n H_i^2$  has a sum of squares as coefficient and  $\mathbb{R}$  is real,  $\sum_{i=1}^n H_i^2$  is a polynomial of even degree  $\leq 2p - 2$ . Hence, the polynomial  $Q$  has odd degree  $\leq p - 2$  and thus has a root  $x$  in  $\mathbb{R}$ . But then  $-1 = \sum_{i=1}^n H_i(x)^2$ , which contradicts the fact that  $\mathbb{R}$  is real. □

*Remark 2.17.* When  $\mathbb{R} = \mathbb{R}$ ,  $a) \Rightarrow b)$  in Theorem 2.11 is nothing but an algebraic proof of the fundamental theorem of algebra. □

**Notation 2.18. [Modulus]** If  $\mathbb{R}$  is real closed, and  $\mathbb{R}[i] = \mathbb{R}[T]/(T^2 + 1)$ , we can identify  $\mathbb{R}[i]$  with  $\mathbb{R}^2$ . For  $z = a + i b \in \mathbb{R}[i]$ ,  $a \in \mathbb{R}$ ,  $b \in \mathbb{R}$ , we define the **conjugate** of  $z$  by  $\bar{z} = a - i b$ . The **modulus** of  $z = a + i b \in \mathbb{R}[i]$  is  $|z| = \sqrt{a^2 + b^2}$ .  $\square$

**Proposition 2.19.** *Let  $\mathbb{R}$  be a real closed field,  $P \in \mathbb{R}[X]$ . The irreducible factors of  $P$  are linear or have the form*

$$(X - c)^2 + d^2 = (X - c - i d)(X - c + i d), d \neq 0$$

with  $c, d \in \mathbb{R}$ .

**Proof:** Use the fact that  $\mathbb{R}[i]$  is algebraically closed by Theorem 2.11 and that the conjugate of a root of  $P$  is a root of  $P$ .  $\square$

**Exercise 2.9.** Prove that, in a real closed field, a second degree polynomial

$$P = a X^2 + b X + c, a \neq 0$$

has a constant non-zero sign if and only if its **discriminant**  $b^2 - 4 a c$  is negative. Hint: the classical computation over the reals is still valid in a real closed field.

Closed, open and semi-open intervals in  $\mathbb{R}$  will be denoted in the usual way:

$$\begin{aligned} (a, b) &= \{x \in \mathbb{R} \mid a < x < b\}, \\ [a, b] &= \{x \in \mathbb{R} \mid a \leq x \leq b\}, \\ (a, b] &= \{x \in \mathbb{R} \mid a < x \leq b\}, \\ (a, +\infty) &= \{x \in \mathbb{R} \mid a < x\}, \\ &\dots \end{aligned}$$

**Proposition 2.20.** *Let  $\mathbb{R}$  be a real closed field,  $P \in \mathbb{R}[X]$  such that  $P$  does not vanish in  $(a, b)$ , then  $P$  has constant sign in the interval  $(a, b)$ .*

**Proof:** Use the fact that  $\mathbb{R}$  has the intermediate value property by Theorem 2.11.  $\square$

This proposition shows that it makes sense to talk about the sign of a polynomial to the right (resp. to the left) of any  $a \in \mathbb{R}$ . Namely, the sign of  $P$  to the right (resp. to the left) of  $a$  is the sign of  $P$  in any interval  $(a, b)$  (resp.  $(b, a)$ ) in which  $P$  does not vanish. We can also speak of the sign of  $P(+\infty)$  (resp.  $P(-\infty)$ ) as the sign of  $P(M)$  for  $M$  sufficiently large (resp. small) i.e. greater (resp. smaller) than any root of  $P$ . This coincides with the sign of  $\text{lcof}(P)$  (resp.  $(-1)^{\deg(P)} \text{lcof}(P)$ ) using Proposition 2.4.

**Proposition 2.21.** *If  $r$  is a root of  $P$  of multiplicity  $\mu$  in a real closed field  $\mathbb{R}$  then the sign of  $P$  to the right of  $r$  is the sign of  $P^{(\mu)}(r)$  and the sign of  $P$  to the left of  $r$  is the sign of  $(-1)^\mu P^{(\mu)}(r)$ .*

**Proof:** Write  $P = (X - r)^\mu Q(x)$  where  $Q(r) \neq 0$ , and note that

$$\text{sign}(Q(r)) = \text{sign}(P^{(\mu)}(r)).$$

□

We next show that univariate polynomials over a real closed field  $\mathbb{R}$  share some of the well known basic properties possessed by differentiable functions over  $\mathbb{R}$ .

**Proposition 2.22. [Rolle’s theorem]** *Let  $\mathbb{R}$  be a real closed field,  $P \in \mathbb{R}[X]$ ,  $a, b \in \mathbb{R}$  with  $a < b$  and  $P(a) = P(b) = 0$ . Then the derivative polynomial  $P'$  has a root in  $(a, b)$ .*

**Proof:** One may reduce to the case where  $a$  and  $b$  are two consecutive roots of  $P$ , i.e. when  $P$  never vanishes on  $(a, b)$ . Then  $P = (X - a)^m (X - b)^n Q$ , where  $Q$  never vanishes on  $[a, b]$ . Thus  $Q$  has constant sign on  $[a, b]$  by Proposition 2.20. Then  $P' = (X - a)^{m-1} (X - b)^{n-1} Q_1$ , where

$$Q_1 = m(X - b)Q + n(X - a)Q + (X - a)(X - b)Q'.$$

Thus  $Q_1(a) = m(a - b)Q(a)$  and  $Q_1(b) = n(b - a)Q(b)$ , and hence  $Q_1(a)$  and  $Q_1(b)$  have opposite signs. By the intermediate value property,  $Q_1$  has a root in  $(a, b)$ , and so does  $P'$ . □

**Corollary 2.23. [Mean Value theorem]** *Let  $\mathbb{R}$  be a real closed field,  $P \in \mathbb{R}[X]$ ,  $a, b \in \mathbb{R}$  with  $a < b$ . There exists  $c \in (a, b)$  such that*

$$P(b) - P(a) = (b - a)P'(c).$$

**Proof:** Apply Rolle’s theorem (Proposition 2.22) to

$$Q(X) = (P(b) - P(a))(X - a) - (b - a)(P(X) - P(a)).$$
□

**Corollary 2.24.** *Let  $\mathbb{R}$  be a real closed field,  $P \in \mathbb{R}[X]$ ,  $a, b \in \mathbb{R}$  with  $a < b$ . If the derivative polynomial  $P'$  is positive (resp. negative) over  $(a, b)$ , then  $P$  is increasing (resp. decreasing) over  $[a, b]$ .*

The following Proposition 2.28 (Basic Thom’s Lemma) which will have important consequences in Chapter 10. We first need a few definitions.

**Definition 2.25.** Let  $\mathcal{Q}$  be a finite subset of  $\mathbb{R}[X_1, \dots, X_k]$ . A **sign condition** on  $\mathcal{Q}$  is an element of  $\{0, 1, -1\}^{\mathcal{Q}}$ , i.e. a mapping from  $\mathcal{Q}$  to  $\{0, 1, -1\}$ . A **strict sign condition** on  $\mathcal{Q}$  is an element of  $\{1, -1\}^{\mathcal{Q}}$ , i.e. a mapping from  $\mathcal{Q}$  to  $\{1, -1\}$ . We say that  $\mathcal{Q}$  **realizes** the sign condition  $\sigma$  at  $x \in \mathbb{R}^k$  if  $\bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(x)) = \sigma(Q)$ .



The realization of the sign condition  $\sigma$  is

$$\text{Reali}(\sigma) = \{x \in \mathbb{R}^k \mid \bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(x)) = \sigma(Q)\}.$$

The sign condition  $\sigma$  is **realizable** if  $\text{Reali}(\sigma)$  is non-empty.  $\square$

**Notation 2.26. [Derivatives]** Let  $P$  be a univariate polynomial of degree  $p$  in  $\mathbb{R}[X]$ . We denote by  $\text{Der}(P)$  the list  $P, P', \dots, P^{(p)}$ .  $\square$

**Proposition 2.27. [Basic Thom's Lemma]** *Let  $P$  be a univariate polynomial of degree  $p$  and let  $\sigma$  be a sign condition on  $\text{Der}(P)$ . Then  $\text{Reali}(\sigma)$  is either empty, a point, or an open interval.*

**Proof:** The proof is by induction on the degree  $p$  of  $P$ . There is nothing to prove if  $p = 0$ . Suppose that the proposition has been proved for  $p - 1$ . Let  $\sigma \in \{0, 1, -1\}^{\text{Der}(P)}$  be a sign condition on  $\text{Der}(P)$ , and let  $\sigma'$  be its restriction to  $\text{Der}(P')$ . If  $\text{Reali}(\sigma')$  is either a point or empty, then

$$\text{Reali}(\sigma) = \text{Reali}(\sigma') \cap \{x \in \mathbb{R} \mid \text{sign}(P(x)) = \sigma(P)\}$$

is either a point or empty. If  $\text{Reali}(\sigma')$  is an open interval,  $P'$  has a constant non-zero sign on it. Thus  $P$  is strictly monotone on  $\text{Reali}(\sigma')$  so that the claimed properties are satisfied for  $\text{Reali}(\sigma)$ .  $\square$

Proposition 2.27 has interesting consequences. One of them is the fact that a root  $x \in \mathbb{R}$  of a polynomial  $P$  of degree  $d$  with coefficients in  $\mathbb{R}$  may be distinguished from the other roots of  $P$  in  $\mathbb{R}$  by the signs of the derivatives of  $P$  at  $x$ .

**Proposition 2.28. [Thom encoding]** *Let  $P$  be a non-zero polynomial of degree  $d$  with coefficients in  $\mathbb{R}$ . Let  $x$  and  $x'$  be two elements of  $\mathbb{R}$ , and denote by  $\sigma$  and  $\sigma'$  the sign conditions on  $\text{Der}(P)$  realized at  $x$  and  $x'$ . Then:*

- If  $\sigma = \sigma'$  with  $\sigma(P) = \sigma'(P) = 0$  then  $x = x'$ .
- If  $\sigma \neq \sigma'$ , one can decide whether  $x < x'$  or  $x > x'$  as follows. Let  $k$  be the smallest integer such that  $\sigma(P^{(d-k)})$  and  $\sigma'(P^{(d-k)})$  are different. Then
  - $\sigma(P^{(d-k+1)}) = \sigma'(P^{(d-k+1)}) \neq 0$ .
  - If  $\sigma(P^{(d-k+1)}) = \sigma'(P^{(d-k+1)}) = 1$ ,
 
$$x > x' \Leftrightarrow \sigma(P^{(d-k)}) > \sigma'(P^{(d-k)}).$$
  - If  $\sigma(P^{(d-k+1)}) = \sigma'(P^{(d-k+1)}) = -1$ ,
 
$$x > x' \Leftrightarrow \sigma(P^{(d-k)}) < \sigma'(P^{(d-k)}).$$

**Proof:** The first item is a consequence of Proposition 2.27. The first part of the second item follows from Proposition 2.27 applied to  $P^{(d-k+1)}$ . The two last parts follow easily since the set

$$\{x \in \mathbb{R} \mid \text{sign}(P^{(i)}(x)) = \sigma(P^{(i)}), i = d - k + 1, \dots, n - 1\}$$

is an interval by Proposition 2.28 applied to  $P^{(d-k+1)}$ , and, on an interval, the sign of the derivative of a polynomial determines whether it is increasing or decreasing.  $\square$

**Definition 2.29.** Let  $P \in \mathbf{R}[X]$  and  $\sigma \in \{0, 1, -1\}^{\text{Der}(P)}$ , a sign condition on the set  $\text{Der}(P)$  of derivatives of  $P$ . The sign condition  $\sigma$  is a **Thom encoding of  $x \in \mathbf{R}$**  if  $\sigma(P) = 0$  and  $\text{Reali}(\sigma) = \{x\}$ , i.e.  $\sigma$  is the sign condition taken by the set  $\text{Der}(P)$  at  $x$ .  $\square$

*Example 2.30.* In any real closed field  $\mathbf{R}$ ,  $P = X^2 - 2$  has two roots, characterized by the sign of the derivative  $2X$ : one root for which  $2X > 0$  and one root for which  $2X < 0$ . Note that no numerical information about the roots is needed to characterize them this way.  $\square$

Any ordered field can be embedded in a real closed field. More precisely, any ordered field  $\mathbf{F}$  possesses a unique **real closure** which is the smallest real closed field extending it. The elements of the real closure are algebraic over  $\mathbf{F}$  (i.e. satisfy an equation with coefficients in  $\mathbf{F}$ ). We refer the reader to [26] for these results.

**Exercise 2.10.** If  $\mathbf{F}$  is contained in a real closed field  $\mathbf{R}$ , the real closure of  $\mathbf{F}$  consists of the elements of  $\mathbf{R}$  which are algebraic over  $\mathbf{F}$ . (Hint: given  $\alpha$  and  $\beta$  roots of  $P$  and  $Q$  in  $\mathbf{F}[X]$ , find polynomials in  $\mathbf{F}[X]$  with roots  $\alpha + \beta$  and  $\alpha\beta$ , using Proposition 2.16).

**Exercise 2.11.** Prove that  $\mathbb{R}_{\text{alg}}$  is real closed. Prove that the field  $\mathbb{R}_{\text{alg}}$  is the real closure of  $\mathbb{Q}$ .

The following theorem proves that any algebraically closed field of characteristic zero is the algebraic closure of a real closed field.

**Theorem 2.31.** *If  $C$  is an algebraically closed field of characteristic zero, there exists a real closed field  $\mathbf{R} \subset C$  such that  $\mathbf{R}[i] = C$ .*

**Proof:** The field  $C$  contains a real subfield, the field  $\mathbb{Q}$  of rational numbers. Let  $\mathbf{R}$  be a maximal real subfield of  $C$ . The field  $\mathbf{R}$  is real closed since it has no nontrivial real algebraic extension contained in  $C$  (see Theorem 2.11). Note that  $C \setminus \mathbf{R}$  cannot contain a  $t$  which is transcendental over  $\mathbf{R}$  since otherwise  $\mathbf{R}(t)$  would be a real field properly containing  $\mathbf{R}$ .  $\square$

An ordered field  $\mathbf{F}$  is **archimedean** if, whenever  $a, b$  are positive elements of  $\mathbf{F}$ , there exists a natural number  $n \in \mathbb{N}$  so that  $na > b$ .

Real closed fields are not necessarily archimedean and may contain infinitesimal elements. We shall see at the end of this chapter an example of a non-archimedean real closed field when we study the field of Puiseux series.

## 2.2 Real Root Counting

Although we have a very simple criterion for determining whether a polynomial  $P \in \mathbb{C}[X]$  has a root in  $\mathbb{C}$  (namely, if and only if  $\deg(P) \neq 0$ ), it is much more difficult to decide whether a polynomial  $P \in \mathbb{R}[X]$  has a root in  $\mathbb{R}$ . The first result in this direction was found more than 350 years ago by Descartes. We begin the section with a generalization of this result.

### 2.2.1 Descartes's Law of Signs and the Budan-Fourier Theorem

**Notation 2.32. [Sign variations]** The **number of sign variations**,  $\text{Var}(a)$ , in a sequence,  $a = a_0, \dots, a_p$ , of elements in  $\mathbb{R} \setminus \{0\}$  is defined by induction on  $p$  by:

$$\begin{aligned} \text{Var}(a_0) &= 0 \\ \text{Var}(a_0, \dots, a_p) &= \begin{cases} \text{Var}(a_1, \dots, a_p) + 1 & \text{if } a_0 a_1 < 0 \\ \text{Var}(a_1, \dots, a_p) & \text{if } a_0 a_1 > 0 \end{cases} \end{aligned}$$

This definition extends to any finite sequence  $a$  of elements in  $\mathbb{R}$  by considering the finite sequence  $b$  obtained by dropping the zeros in  $a$  and defining

$$\text{Var}(a) = \text{Var}(b), \quad \text{Var}(\emptyset) = 0.$$

For example  $\text{Var}(1, -1, 2, 0, 0, 3, 4, -5, -2, 0, 3) = 4$ . □

Let  $P = a_p X^p + \dots + a_0$  be a univariate polynomial in  $\mathbb{R}[X]$ . We write  $\text{Var}(P)$  for the number of sign variations in  $a_0, \dots, a_p$  and  $\text{pos}(P)$  for the number of positive real roots of  $P$ , counted with multiplicity.

### Theorem 2.33. [Descartes' law of signs]

- $\text{Var}(P) \geq \text{pos}(P)$
- $\text{Var}(P) - \text{pos}(P)$  is even.

We will prove the following generalization of Theorem 2.33 (Descartes's law of signs) due to Budan and Fourier.

### Notation 2.34. [Sign variations in a sequence of polynomials at $a$ ]

Let  $\mathcal{P} = P_0, P_1, \dots, P_d$  be a sequence of polynomials and let  $a$  be an element of  $\mathbb{R} \cup \{-\infty, +\infty\}$ . The **number of sign variations** of  $\mathcal{P}$  at  $a$ , denoted by  $\text{Var}(\mathcal{P}; a)$ , is  $\text{Var}(P_0(a), \dots, P_d(a))$  (at  $-\infty$  and  $+\infty$  the signs to consider are the signs of the leading monomials according to Proposition 2.4).

For example, if  $\mathcal{P} = X^5, X^2 - 1, 0, X^2 - 1, X + 2, 1$ ,  $\text{Var}(\mathcal{P}; 1) = 0$ .

Given  $a$  and  $b$  in  $\mathbb{R} \cup \{-\infty, +\infty\}$ , we denote

$$\text{Var}(\mathcal{P}; a, b) = \text{Var}(\mathcal{P}; a) - \text{Var}(\mathcal{P}; b).$$

□

We denote by  $\text{num}(P; (a, b])$  the number of roots of  $P$  in  $(a, b]$  counted with multiplicities.

**Theorem 2.35. [Budan-Fourier theorem]** *Let  $P$  be a univariate polynomial of degree  $p$  in  $\mathbb{R}[X]$ . Given  $a$  and  $b$  in  $\mathbb{R} \cup \{-\infty, +\infty\}$*

- $\text{Var}(\text{Der}(P); a, b) \geq \text{num}(P; (a, b])$ ,
- $\text{Var}(\text{Der}(P); a, b) - \text{num}(P; (a, b])$  is even.

Theorem 2.33 (Descartes’s law of signs) is a particular case of Theorem 2.35 (Budan-Fourier).

**Proof of Theorem 2.33 (Descartes’ law of signs):** The coefficient of degree  $i$  of  $P$  has the same sign as the  $p - i$ -th derivative of  $P$  evaluated at 0. Moreover, there are no sign variations in the signs of the derivatives at  $+\infty$ . So that  $\text{Var}(P) = \text{Var}(\text{Der}(P); 0, +\infty)$ .  $\square$

The following lemma is the key to the proof of Theorem 2.35 (Budan-Fourier).

**Lemma 2.36.** *Let  $c$  be a root of  $P$  of multiplicity  $\mu \geq 0$ . If no  $P^{(k)}$ ,  $0 \leq k \leq p$ , has a root in  $[d, c) \cup (c, d']$ , then*

- a)  $\text{Var}(\text{Der}(P); d, c) - \mu$  is non-negative and even,
- b)  $\text{Var}(\text{Der}(P); c, d') = 0$ .

**Proof:** We prove the claim by induction on the degree of  $P$ . The claim is true if the degree of  $P$  is 1.

Suppose first that  $P(c) = 0$ , and hence  $\mu > 0$ . By induction hypothesis applied to  $P'$ ,

- a)  $\text{Var}(\text{Der}(P'); d, c) - (\mu - 1)$  is non-negative and even,
- b)  $\text{Var}(\text{Der}(P'); c, d') = 0$ .

The sign of  $P$  at the left of  $c$  is the opposite of the sign of  $P'$  at the left of  $c$  and the sign of  $P$  at the right of  $c$  is the sign of  $P'$  at the right of  $c$ . Thus

$$\begin{aligned} \text{Var}(\text{Der}(P); d) &= \text{Var}(\text{Der}(P'); d) + 1, \\ \text{Var}(\text{Der}(P); c) &= \text{Var}(\text{Der}(P'); c), \\ \text{Var}(\text{Der}(P); d') &= \text{Var}(\text{Der}(P'); d'), \end{aligned} \tag{2.1}$$

and the claim follows.

Suppose now that  $P(c) \neq 0$ , and hence  $\mu = 0$ . Let  $\nu$  be the multiplicity of  $c$  as a root of  $P'$ . By induction hypothesis applied to  $P'$

- a)  $\text{Var}(\text{Der}(P'); d, c) - \nu$  is non-negative and even,
- b)  $\text{Var}(\text{Der}(P'); c, d') = 0$ .

There are four cases to consider.

If  $\nu$  is odd, and  $\text{sign}(P^{(\nu+1)}(c)P(c)) > 0$ ,

$$\begin{aligned}\text{Var}(\text{Der}(P); d) &= \text{Var}(\text{Der}(P'); d) + 1, \\ \text{Var}(\text{Der}(P); c) &= \text{Var}(\text{Der}(P'); c), \\ \text{Var}(\text{Der}(P); d') &= \text{Var}(\text{Der}(P'); d').\end{aligned}\tag{2.2}$$

If  $\nu$  is odd, and  $\text{sign}(P^{(\nu+1)}(c)P(c)) < 0$ ,

$$\begin{aligned}\text{Var}(\text{Der}(P); d) &= \text{Var}(\text{Der}(P'); d), \\ \text{Var}(\text{Der}(P); c) &= \text{Var}(\text{Der}(P'); c) + 1, \\ \text{Var}(\text{Der}(P); d') &= \text{Var}(\text{Der}(P'); d') + 1.\end{aligned}\tag{2.3}$$

If  $\nu$  is even, and  $\text{sign}(P^{(\nu+1)}(c)P(c)) > 0$ ,

$$\begin{aligned}\text{Var}(\text{Der}(P); d) &= \text{Var}(\text{Der}(P'); d), \\ \text{Var}(\text{Der}(P); c) &= \text{Var}(\text{Der}(P'); c), \\ \text{Var}(\text{Der}(P); d') &= \text{Var}(\text{Der}(P'); d').\end{aligned}\tag{2.4}$$

If  $\nu$  is even, and  $\text{sign}(P^{(\nu+1)}(c)P(c)) < 0$ ,

$$\begin{aligned}\text{Var}(\text{Der}(P); d) &= \text{Var}(\text{Der}(P'); d) + 1, \\ \text{Var}(\text{Der}(P); c) &= \text{Var}(\text{Der}(P'); c) + 1, \\ \text{Var}(\text{Der}(P); d') &= \text{Var}(\text{Der}(P'); d') + 1.\end{aligned}\tag{2.5}$$

The claim is true in each of these four cases.  $\square$

**Proof of Theorem 2.35:** It is clear that, for every  $c \in (a, b)$ ,

$$\begin{aligned}\text{num}(P; (a, b]) &= \text{num}(P; (a, c]) + \text{num}(P; (c, b]) \\ \text{Var}(\text{Der}(P); a, b) &= \text{Var}(\text{Der}(P); a, c) + \text{Var}(\text{Der}(P); c, b).\end{aligned}$$

Let  $c_1 < \dots < c_r$  be the roots of all the polynomials  $P^{(j)}$ ,  $0 \leq j \leq p-1$ , in the interval  $(a, b)$  and let  $a = c_0, b = c_{r+1}, d_i \in (c_i, c_{i+1})$  so that

$$a = c_0 < d_0 < c_1 < \dots < c_r < d_r < c_{r+1} = b.$$

Since,

$$\begin{aligned}\text{num}(P; (a, b]) &= \sum_{i=0}^r \text{num}(P; (c_i, d_i]) + \text{num}(P; (d_i, c_{i+1}]), \\ \text{Var}(\text{Der}(P); a, b) &= \sum_{i=0}^r \text{Var}(\text{Der}(P); c_i, d_i) + \text{Var}(\text{Der}(P); d_i, c_{i+1}),\end{aligned}$$

the claim follows immediately from Lemma 2.36.  $\square$

In general it is not possible to conclude much about the number of roots on an interval using only Theorem 2.35 (Descartes's law of signs).

*Example 2.37.* The polynomial  $P = X^2 - X + 1$  has no real root, but  $\text{Var}(\text{Der}(P); 0, 1) = 2$ . It is impossible to find  $a \in (0, 1]$  such that  $\text{Var}(\text{Der}(P); 0, a) = 1$  and  $\text{Var}(\text{Der}(P); a, 1) = 1$  since otherwise  $P$  would have two real roots. This means that however we refine the interval  $(0, 1]$ , we are going to have an interval (the interval  $(a, b]$  containing  $1/2$ ) giving 2 sign variations.  $\square$

However, there are particular cases where Theorem 2.35 (Budan-Fourier) gives the number of roots on an interval:

**Exercise 2.12.** Prove that

- If  $\text{Var}(\text{Der}(P); a, b) = 0$ , then  $P$  has no root in  $(a, b]$ .
- If  $\text{Var}(\text{Der}(P); a, b) = 1$ , then  $P$  has exactly one root in  $(a, b]$ , which is simple.

*Remark 2.38.* Another important instance, used in Chapter 8, where Theorem 2.35 (Budan-Fourier) permits a sharp conclusion is the following. When we know in advance that all the roots of a polynomial are real, i.e. when  $\text{num}(P; (-\infty, +\infty)) = p$ , the number  $\text{Var}(\text{Der}(P); a, b)$  is exactly the number of roots counted with multiplicities in  $(a, b]$ . Indeed the number  $\text{Var}(\text{Der}(P); -\infty, +\infty)$ , which is always at most  $p$ , is here equal to  $p$ , hence

$$\begin{aligned} \text{num}(P; (-\infty, a]) &\leq \text{Var}(\text{Der}(P); -\infty, a) \\ \text{num}(P; (a, b]) &\leq \text{Var}(\text{Der}(P); a, b) \\ \text{num}(P; (b, +\infty)) &\leq \text{Var}(\text{Der}(P); b, +\infty) \end{aligned}$$

imply  $\text{num}(P, (a, b]) = \text{Var}(\text{Der}(P); a, b)$ .  $\square$

We are going now to describe situations where the number of sign variations in the coefficients coincides exactly with the number of real roots.

The first case we consider is obvious.

**Proposition 2.39.** *Let  $P \in \mathbb{R}[X]$  be a monic polynomial. If all the roots of  $P$  have non-positive real part, then  $\text{Var}(P) = 0$ .*

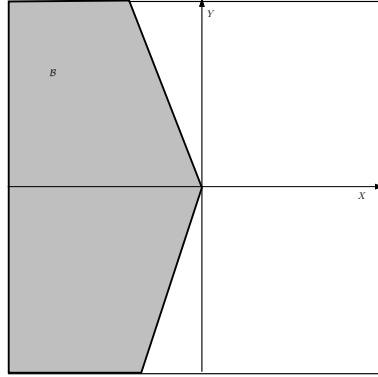
**Proof:** Obvious, using the decomposition of  $P$  in products of linear factors and polynomials of degree 2 with complex conjugate roots, since the product of two polynomials whose coefficients are all non-negative have coefficients that are all non-negative.  $\square$

The second case we consider is the case of normal polynomials. A polynomial  $A = a_p X^p + \dots + a_0$  with non-negative coefficients is **normal** if

- a)  $a_p > 0$ ,
- b)  $a_k^2 \geq a_{k-1} a_{k+1}$  for all index  $k$ ,
- c)  $a_h > 0$  and  $a_j > 0$  for indices  $j < h$  implies  $a_{j+1} > 0, \dots, a_{h-1} > 0$

(with the convention that  $a_i = 0$  if  $i < 0$  or  $i > p$ ).

**Proposition 2.40.** *Let  $P \in \mathbb{R}[X]$  be a monic polynomial. If all the roots of  $P$  belong to the cone  $\mathcal{B}$  of the complex plane (see Figure 2.1) defined by  $\mathcal{B} = \{a + ib \mid |b| \leq -\sqrt{3}a\}$ , then  $P$  is normal.*



**Fig. 2.1.** Cone  $\mathcal{B}$

The proof of Proposition 2.40 relies on the following lemmas.

**Lemma 2.41.** *The polynomial  $X - x$  is normal if and only if  $x \leq 0$ .*

**Proof:** Follows immediately from the definition of a normal polynomial.  $\square$

**Lemma 2.42.** *A quadratic monic polynomial  $A$  with complex conjugate roots is normal if and only if its roots belong to the cone  $\mathcal{B}$ .*

**Proof:**

Let  $a + ib$  and  $a - ib$  be the roots of  $A$ . Then

$$A = X^2 - 2aX + (a^2 + b^2)$$

is normal if and only if

- a)  $-2a \geq 0$ ,
- b)  $a^2 + b^2 \geq 0$ ,
- c)  $(-2a)^2 \geq a^2 + b^2$ .

that is if and only if  $a \leq 0$  and  $4a^2 \geq a^2 + b^2$ , or equivalently  $a + ib \in \mathcal{B}$ .  $\square$

**Lemma 2.43.** *The product of two normal polynomials is normal.*

**Proof:** Let  $A = a_p X^p + \dots + a_0$  and  $B = b_q X^q + \dots + b_0$  be two normal polynomials. We can suppose without loss of generality that 0 is not a root of  $A$  and  $B$ , i.e. that all the coefficients of  $A$  and  $B$  are positive.

Let  $C = AB = c_{p+q} X^{p+q} + \dots + c_0$ . It is clear that all the coefficients of  $C$  are positive.

It remains to prove that  $c_k^2 \geq c_{k-1} c_{k+1}$ .

Using the partition of  $\{(h, j) \in \mathbb{Z}^2 \mid h > j\}$  in  $\{(j+1, h-1) \in \mathbb{Z}^2 \mid h \leq j\}$  and  $\{(h, h-1) \mid h \in \mathbb{Z}\}$ .

$$\begin{aligned}
 c_k^2 - c_{k-1} c_{k+1} &= \sum_{h \leq j} a_h a_j b_{k-h} b_{k-j} + \sum_{h > j} a_h a_j b_{k-h} b_{k-j} \\
 &\quad - \sum_{h \leq j} a_h a_j b_{k-h+1} b_{k-j-1} - \sum_{h > j} a_h a_j b_{k-h+1} b_{k-j-1} \\
 &= \sum_{h \leq j} a_h a_j b_{k-h} b_{k-j} + \sum_{h \leq j} a_{j+1} a_{h-1} b_{k-j-1} b_{k-h-1} \\
 &\quad + \sum_h a_h a_{h-1} b_{k-h} b_{k-h+1} - \sum_h a_h a_{h-1} b_{k-h+1} b_{k-h} \\
 &\quad - \sum_{h \leq j} a_h a_j b_{k-h+1} b_{k-j-1} - \sum_{h \leq j} a_{j+1} a_{h-1} b_{k-j} b_{k-h} \\
 &= \sum_{h \leq j} (a_h a_j - a_{h-1} a_{j+1}) (b_{k-j} b_{k-h} - b_{k-j-1} b_{k-h+1}).
 \end{aligned}$$

Since  $A$  is normal and  $a_0, \dots, a_p$  are positive, one has

$$\frac{a_{p-1}}{a_p} \geq \frac{a_{p-2}}{a_{p-1}} \geq \dots \geq \frac{a_0}{a_1},$$

and  $a_h a_j - a_{h-1} a_{j+1} \geq 0$ , for all  $k \leq j$ . Similar inequalities hold for the coefficients of  $B$  and finally  $c_k^2 - c_{k-1} c_{k+1}$  is non-negative, being a sum of non-negative quantities. □

**Proof of Proposition 2.40:** Factor  $P$  into linear and quadratic polynomials. By Lemma 2.41 and Lemma 2.42 each of these factors is normal. Now use Lemma 2.43. □

Finally we obtain the following partial reciprocal to Descartes law of signs.

**Proposition 2.44.** *If  $A$  is normal and  $x > 0$ , then  $\text{Var}(A(X-x)) = 1$ .*

**Proof:** We can suppose without loss of generality that that 0 is not a root of  $A$ , that it that all the coefficients of  $A$  are positive.

Then

$$\frac{a_{p-1}}{a_p} \geq \frac{a_{p-2}}{a_{p-1}} \geq \dots \geq \frac{a_0}{a_1},$$

and

$$\frac{a_{p-1}}{a_p} - x \geq \frac{a_{p-2}}{a_{p-1}} - x \geq \dots \geq \frac{a_0}{a_1} - x.$$

Since  $a_p > 0$  and  $-a_0 x < 0$ , the coefficients of the polynomial

$$(X-x)A = a_p X^{p+1} + a_p \left( \frac{a_{p-1}}{a_p} - x \right) X^p + \dots + a_1 \left( \frac{a_0}{a_1} - x \right) X - a_0 x.$$

have exactly one sign variation. □

A natural question when looking at Budan-Fourier's Theorem (Theorem 2.35), is to interpret the even difference  $\text{Var}(\text{Der}(P); a, b) - \text{num}(P; (a, b])$ . This can be done through the notion of virtual roots.



The virtual roots of  $P$  will enjoy the following properties:

- a) the number of virtual roots of  $P$  counted with virtual multiplicities is equal to the degree  $p$  of  $P$ ,
- b) on an open interval defined by virtual roots, the sign of  $P$  is fixed,
- c) virtual roots of  $P$  and virtual roots of  $P'$  are interlaced: if  $x_1 \leq \dots \leq x_p$  are the virtual roots of  $P$  and  $y_1 \leq \dots \leq y_{p-1}$  are the virtual roots of  $P'$ , then

$$x_1 \leq y_1 \leq \dots \leq x_{p-1} \leq y_{p-1} \leq x_p.$$

Given these properties, in the particular case where  $P$  is a polynomial of degree  $p$  with all its roots real and simple, virtual roots and real roots clearly coincide.

**Definition 2.45. [Virtual roots]** The definition of **virtual roots** proceeds by induction on  $p = \deg(P)$ . We prove simultaneously that properties a), b), c) hold.

If  $p = 0$ ,  $P$  has no virtual root and properties a), b), c) hold.

Suppose that properties a), b), c) hold for the virtual roots of  $P'$ .

By induction hypothesis the virtual roots of  $P'$  are  $y_1 \leq \dots \leq y_{p-1}$ . Let

$$I_1 = (-\infty, y_1], \dots, I_i = [y_{i-1}, y_i], \dots, I_p = [y_{p-1}, +\infty).$$

By induction hypothesis, the sign of  $P'$  is fixed on the interior of each  $I_i$ .

Let  $x_i$  be unique value in  $I_i$  such that the absolute value of  $P$  on  $I_i$  reaches its minimum. The virtual roots of  $P$  are  $x_1 \leq \dots \leq x_p$ .

According to this inductive definition, properties a), b) and c) are clear for virtual roots of  $P$ . Note that the virtual roots of  $P$  are always roots of a derivative of  $P$ .

The **virtual multiplicity** of  $x$  with respect to  $P$ , denoted  $v(P, x)$  is the number of times  $x$  is repeated in the list  $x_1 \leq \dots \leq x_p$  of virtual roots of  $P$ . In particular, if  $x$  is not a virtual root of  $P$ , its virtual multiplicity is equal to 0. Note that if  $x$  is a virtual root of  $P'$  with virtual multiplicity  $\nu$  with respect to  $P$ , the virtual multiplicity of  $x$  with respect to  $P'$  can only be  $\nu$ ,  $\nu + 1$  or  $\nu - 1$ . Moreover, if  $x$  is a root of  $P'$ , the virtual multiplicity of  $x$  with respect to  $P'$  is necessarily  $\nu + 1$ .  $\square$

*Example 2.46.* The virtual roots of a polynomial  $P$  of degree 2 are

- the two roots of  $P$  with virtual multiplicity 1 if  $P$  has two distinct real roots,
- the root of  $P'$  with virtual multiplicity 2 if  $P$  does not have two distinct real roots.  $\square$

Given  $a$  and  $b$ , we denote by  $v(P; (a, b])$  the number of virtual roots of  $P$  in  $(a, b]$  counted with virtual multiplicities.

**Theorem 2.47.**

$$v(P; (a, b]) = \text{Var}(\text{Der}(P); a, b).$$

The following lemma is the key to the proof of Theorem 2.47.

**Lemma 2.48.** *Let  $c$  be a root of  $P$  of virtual multiplicity  $v(P, c) \geq 0$ . If no  $P^{(k)}$ ,  $0 \leq k < p$  has a root in  $[d, c)$ , then*

$$v(P, c) = \text{Var}(\text{Der}(P); d, c).$$

**Proof:** The proof of the claim is by induction on  $p = \deg(P)$ . The claim obviously holds if  $p = 0$ .

Let  $w = v(P, c)$ .

- If  $c$  is a root of  $P$ , the virtual multiplicity of  $c$  as a root of  $P'$  is  $w - 1$ . By induction hypothesis applied to  $P'$ ,  $\text{Var}(\text{Der}(P'); d, c) = w - 1$ . The claim follows from equation (2.1).
- If  $c$  is not a root of  $P$ , is a virtual root of  $P$  with virtual multiplicity  $w$ , and a virtual root of  $P'$  with multiplicity  $\nu$  and virtual multiplicity  $u$ , by induction hypothesis applied to  $P'$ ,  $\text{Var}(\text{Der}(P'); d, c) = u$ .
  - If the sign of  $P'$  at the left and at the right of  $c$  differ,  $\nu$  is odd as well as  $u$ , using Lemma 2.36 a) and the induction hypothesis for  $P'$ .
    - If  $c$  is a local minimum of the absolute value of  $P$ ,  $w = u + 1$ ,  $\text{sign}(P^{(\nu+1)}(c) P(c)) > 0$ , and the claim follows from (2.2).
    - If  $c$  is a local maximum of the absolute value of  $P$ ,  $w = u - 1$ ,  $\text{sign}(P^{(\nu+1)}(c) P(c)) < 0$ , and the claim follows from (2.3).
  - If the sign of  $P'$  at the left and at the right of  $c$  coincide,  $w = u$ ,  $\nu$  is even as well as  $u$  using Lemma 2.36 a) and the induction hypothesis for  $P'$ . The claim follows from (2.4) and (2.5).

The claim follows in each of these cases. □

It follows clearly from Proposition 2.48 that:

**Corollary 2.49.** *All the roots of  $P$  are virtual roots of  $P$ . The virtual multiplicity is at least equal to the multiplicity and the difference is even.*

**Proof of Theorem 2.47:** It is clear that, for every  $c \in (a, b)$ ,

$$\begin{aligned} v(P; (a, b)) &= v(P; (a, c]) + v(P; (c, b]), \\ \text{Var}(\text{Der}(P); a, b) &= \text{Var}(\text{Der}(P); a, c) + \text{Var}(\text{Der}(P); c, b). \end{aligned}$$

Let  $c_1 < \dots < c_r$  be the roots of all the  $P^{(i)}$ ,  $0 \leq i \leq p - 1$ , in the interval  $(a, b)$  and let  $c_0 = \infty, c_{r+1} = +\infty, d_i \in (c_i, c_{i+1})$  so that  $c_0 < d_0 < c_1 < \dots < c_r < d_r < c_{r+1}$ .

Since

$$\begin{aligned} v(P; (a, b)) &= \sum_{i=0}^r (v(P; (c_i, d_i]) + v(P; (d_i, c_{i+1}))), \\ \text{Var}(\text{Der}(P); a, b) &= \sum_{i=0}^r (\text{Var}(\text{Der}(P); c_i, d_i) + \text{Var}(\text{Der}(P); d_i, c_{i+1})), \end{aligned}$$

the claim follows immediately from Lemma 2.36 b) and Lemma 2.48. □

Finally the even number  $\text{Var}(\text{Der}(P); a, b) \geq \text{num}(P; (a, b])$  appearing in the statement of Budan-Fourier's Theorem (Theorem 2.35) is the sum of the differences between virtual multiplicities and multiplicities of roots of  $P$  in  $(a, b]$ .

### 2.2.2 Sturm's Theorem and the Cauchy Index

Let  $P$  be a non-zero polynomial with coefficients in a real closed field  $\mathbb{R}$ . The sequence of signed remainders of  $P$  and  $P'$ ,  $\text{SRemS}(P, P')$  (see Definition 1.7) is the **Sturm sequence** of  $P$ .

We will prove that the number of roots of  $P$  in  $(a, b)$  can be computed from the Sturm sequence  $\text{SRemS}(P, P')$  evaluated at  $a$  and  $b$  (see Notation 2.34). More precisely the number of roots of  $P$  in  $(a, b)$  is the difference in the number of sign variations in the Sturm's sequence  $\text{SRemS}(P, P')$  evaluated at  $a$  and  $b$ .

**Theorem 2.50. [Sturm's theorem]** *Given  $a$  and  $b$  in  $\mathbb{R} \cup \{-\infty, +\infty\}$ ,*

$$\text{Var}(\text{SRemS}(P, P'); a, b)$$

*is the number of roots of  $P$  in the interval  $(a, b)$ .*

*Remark 2.51.* As a consequence, we can decide whether  $P$  has a root in  $\mathbb{R}$  by checking whether  $\text{Var}(\text{SRemS}(P, P'); -\infty, +\infty) > 0$ .  $\square$

Let us first see how to use Theorem 2.50 (Sturm's theorem).

*Example 2.52.* Consider the polynomial  $P = X^4 - 5X^2 + 4$ . The Sturm sequence of  $P$  is

$$\begin{aligned} \text{SRemS}_0(P, P') &= P = X^4 - 5X^2 + 4, \\ \text{SRemS}_1(P, P') &= P' = 4X^3 - 10X, \\ \text{SRemS}_2(P, P') &= \frac{5}{2}X^2 - 4, \\ \text{SRemS}_3(P, P') &= \frac{18}{5}X, \\ \text{SRemS}_4(P, P') &= 4. \end{aligned}$$

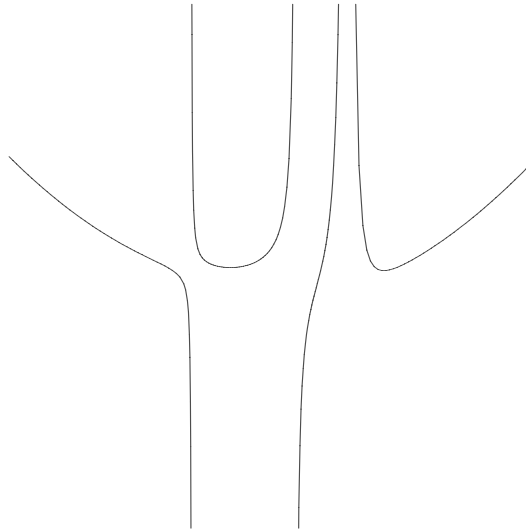
The signs of the leading coefficients of the Sturm sequence are  $++++$  and the degrees of the polynomials in the Sturm sequence are  $4, 3, 2, 1, 0$ . The signs of the polynomials in the Sturm sequence at  $-\infty$  are  $+ - + - +$ , and the signs of the polynomials in the Sturm sequence at  $+\infty$  are  $++++$ , so  $\text{Var}(\text{SRemS}(P, P'); -\infty, +\infty) = 4$ . There are indeed 4 real roots:  $1, -1, 2$ , and  $-2$ .  $\square$

We are going to prove a statement more general than Theorem 2.50 (Sturm's theorem), since it will be useful not only to determine whether  $P$  has a root in  $\mathbb{R}$  but also to determine whether  $P$  has a root at which another polynomial  $Q$  is positive.

With this goal in mind, it is profitable to look at the jumps (discontinuities) of the rational function  $P'Q/P$ . Clearly, these occur only at points  $c$  for which  $P(c) = 0$ ,  $Q(c) \neq 0$ . If  $c$  occurs as a root of  $P$  with multiplicity  $\mu$  then  $P'Q/P = \mu Q(c)/(X - c) + R_c$ , where  $R_c$  is a rational function defined at  $c$ . It is now obvious that if  $Q(c) > 0$ , then  $P'Q/P$  jumps from  $-\infty$  to  $+\infty$  at  $c$ , and if  $Q(c) < 0$ , then  $P'Q/P$  jumps from  $+\infty$  to  $-\infty$  at  $c$ . Thus the number of jumps of  $P'Q/P$  from  $-\infty$  to  $+\infty$  minus the number of jumps of  $P'Q/P$  from  $+\infty$  to  $-\infty$  is equal to the number of roots of  $P$  at which  $Q$  is positive minus the number of roots of  $P$  at which  $Q$  is negative. This observation leads us to the following definition. We need first what we mean by a jump from  $-\infty$  to  $+\infty$ .

**Definition 2.53. [Cauchy index]** Let  $x$  be a root of  $P$ . The function  $Q/P$  **jumps from  $-\infty$  to  $+\infty$  at  $x$**  if the multiplicity  $\mu$  of  $x$  as a root of  $P$  is bigger than the multiplicity  $\nu$  of  $x$  as a root of  $Q$ ,  $\mu - \nu$  is odd and the sign of  $Q/P$  at the right of  $x$  is positive. Similarly, the function  $Q/P$  **jumps from  $+\infty$  to  $-\infty$  at  $x$**  if the multiplicity  $\mu$  of  $x$  as a root of  $P$  is bigger than the multiplicity  $\nu$  of  $x$  as a root of  $Q$ ,  $\mu - \nu$  is odd and the sign of  $Q/P$  at the right of  $x$  is negative.

Given  $a < b$  in  $\mathbb{R} \cup \{-\infty, +\infty\}$  and  $P, Q \in \mathbb{R}[X]$ , we define the **Cauchy index** of  $Q/P$  on  $(a, b)$ ,  $\text{Ind}(Q/P; a, b)$ , to be the number of jumps of the function  $Q/P$  from  $-\infty$  to  $+\infty$  minus the number of jumps of the function  $Q/P$  from  $+\infty$  to  $-\infty$  on the open interval  $(a, b)$ . The **Cauchy index** of  $Q/P$  on  $\mathbb{R}$  is simply called the **Cauchy index** of  $Q/P$  and it is denoted by  $\text{Ind}(Q/P)$ , rather than by  $\text{Ind}(Q/P; -\infty, +\infty)$ .  $\square$



**Fig. 2.2.** Graph of the rational function  $Q/P$

*Example 2.54.* Let

$$\begin{aligned} P &= (X - 3)^2(X - 1)(X + 3), \\ Q &= (X - 5)(X - 4)(X - 2)(X + 1)(X + 2)(X + 4). \end{aligned}$$

The graph of  $Q/P$  is depicted in Figure 2.2.

In this example,

$$\begin{aligned} \text{Ind}(Q/P) &= 0 \\ \text{Ind}(Q/P; -\infty, 0) &= 1 \\ \text{Ind}(Q/P; 0, \infty) &= -1 \end{aligned}$$

□

*Remark 2.55.*

- a) Suppose  $\deg(P) = p$ ,  $\deg(Q) = q < p$ . The Cauchy index  $\text{Ind}(Q/P; a, b)$  is equal to  $p$  if and only if  $q = p - 1$ , the signs of the leading coefficients of  $P$  and  $Q$  are equal, all the roots of  $P$  and  $Q$  are simple and belong to  $(a, b)$ , and there is exactly one root of  $Q$  between two roots of  $P$ .
- b) If  $R = \text{Rem}(Q, P)$ , it follows clearly from the definition that

$$\text{Ind}(Q/P; a, b) = \text{Ind}(R/P; a, b). \quad \square$$

Using the notion of Cauchy index we can reformulate our preceding discussion, using the following notation.

**Notation 2.56. [Tarski-query]** Let  $P \neq 0$  and  $Q$  be elements of  $\mathbb{K}[X]$ . The **Tarski-query** of  $Q$  for  $P$  in  $(a, b)$  is the number

$$\text{TaQ}(Q, P; a, b) = \sum_{x \in (a, b), P(x)=0} \text{sign}(Q(x)).$$

Note that  $\text{TaQ}(Q, P; a, b)$  is equal to

$$\#\{x \in (a, b) \mid P(x) = 0 \wedge Q(x) > 0\} - \#\{x \in (a, b) \mid P(x) = 0 \wedge Q(x) < 0\}$$

where  $\#(S)$  is the number of elements in the finite set  $S$ .

The Tarski-query of  $Q$  for  $P$  on  $\mathbb{R}$  is simply called the **Tarski-query** of  $Q$  for  $P$ , and is denoted by  $\text{TaQ}(Q, P)$ , rather than by  $\text{TaQ}(Q, P; -\infty, +\infty)$ . □

The preceding discussion implies:

**Proposition 2.57.**

$$\text{TaQ}(Q, P; a, b) = \text{Ind}(P'Q/P; a, b).$$

*In particular the number of roots of  $P$  in  $(a, b)$  is  $\text{Ind}(P'/P; a, b)$ .*

We now describe how to compute  $\text{Ind}(Q/P; a, b)$ . We will see that the Cauchy index is the difference in the number of sign variations in the signed remainder sequence  $\text{SRemS}(P, Q)$  evaluated at  $a$  and  $b$  (Definition 1.7 and Notation 2.34).

**Theorem 2.58.** *Let  $P, P \neq 0$ , and  $Q$  be two polynomials with coefficients in a real closed field  $\mathbb{R}$ , and let  $a$  and  $b$  (with  $a < b$ ) be elements of  $\mathbb{R} \cup \{-\infty, +\infty\}$  that are not roots of  $P$ . Then,*

$$\text{Var}(\text{SRemS}(P, Q); a, b) = \text{Ind}(Q/P; a, b).$$

Let  $R = \text{Rem}(P, Q)$  and let  $\sigma(a)$  be the sign of  $PQ$  at  $a$  and  $\sigma(b)$  be the sign of  $PQ$  at  $b$ . The proof of Theorem 2.58 proceeds by induction on the length of the signed remainder sequence and is based on the following lemmas.

**Lemma 2.59.** *If  $a$  and  $b$  are not roots of a polynomial in the signed remainder sequence,*

$$\begin{aligned} & \text{Var}(\text{SRemS}(P, Q); a, b) \\ = & \begin{cases} \text{Var}(\text{SRemS}(Q, -R); a, b) + \sigma(b) & \text{if } \sigma(a)\sigma(b) = -1, \\ \text{Var}(\text{SRemS}(Q, -R); a, b) & \text{if } \sigma(a)\sigma(b) = 1. \end{cases} \end{aligned}$$

**Proof:** The claim follows from the fact that at any  $x$  which is not a root of  $P$  and  $Q$  (and in particular at  $a$  and  $b$ )

$$\text{Var}(\text{SRemS}(P, Q); x) = \begin{cases} \text{Var}(\text{SRemS}(Q, -R); x) + 1 & \text{if } P(x)Q(x) < 0, \\ \text{Var}(\text{SRemS}(Q, -R); x) & \text{if } P(x)Q(x) > 0, \end{cases}$$

looking at all possible cases. □

**Lemma 2.60.** *If  $a$  and  $b$  are not roots of a polynomial in the signed remainder sequence,*

$$\text{Ind}(Q/P; a, b) = \begin{cases} \text{Ind}(-R/Q; a, b) + \sigma(b) & \text{if } \sigma(a)\sigma(b) = -1, \\ \text{Ind}(-R/Q; a, b) & \text{if } \sigma(a)\sigma(b) = 1. \end{cases}$$

**Proof:** We can suppose without loss of generality that  $Q$  and  $P$  are coprime. Indeed if  $D$  is a greatest common divisor of  $P$  and  $Q$  and

$$P_1 = P/D, Q_1 = Q/D, R_1 = \text{Rem}(P_1, Q_1) = R/D,$$

then  $P_1$  and  $Q_1$  are coprime,

$$\text{Ind}(Q/P; a, b) = \text{Ind}(Q_1/P_1; a, b), \text{Ind}(-R/Q; a, b) = \text{Ind}(-R_1/Q_1; a, b),$$

and the signs of  $P(x)Q(x)$  and  $P_1(x)Q_1(x)$  coincide at any point which is not a root of  $PQ$ .

Let  $n_{-+}$  (resp.  $n_{+-}$ ) denote the number of sign variations from  $-1$  to  $1$  (resp. from  $1$  to  $-1$ ) of  $PQ$  when  $x$  varies from  $a$  to  $b$ . It is clear that

$$n_{-+} - n_{+-} = \begin{cases} \sigma(b) & \text{if } \sigma(a)\sigma(b) = -1 \\ 0 & \text{if } \sigma(a)\sigma(b) = 1. \end{cases}$$

It follows from the definition of Cauchy index that

$$\text{Ind}(Q/P; a, b) + \text{Ind}(P/Q; a, b) = n_{-+} - n_{+-}.$$

Noting that

$$\text{Ind}(R/Q; a, b) = \text{Ind}(P/Q; a, b),$$

the claim of the lemma is now clear.  $\square$

**Proof of Theorem 2.58:** We can assume without loss of generality that  $a$  and  $b$  are not roots of a polynomial in the signed remainder sequence. Indeed if  $a < a' < b' < b$  with  $(a, a']$  and  $[b', b)$  containing no root of the polynomials in the signed remainder sequence, it is clear that

$$\text{Ind}(Q/P; a, b) = \text{Ind}(Q/P; a', b').$$

We prove now that

$$\text{Var}(\text{SRemS}(P, Q); a, b) = \text{Var}(\text{SRemS}(P, Q); a', b').$$

We omit  $(P, Q)$  in the notation in the following lines. First notice that since  $a$  is not a root of  $P$ ,  $a$  is not a root of the greatest common divisor of  $P$  and  $Q$ , and hence  $a$  is not simultaneously a root of  $\text{SRemS}_j$  and  $\text{SRemS}_{j+1}$  (resp.  $\text{SRemS}_{j-1}$  and  $\text{SRemS}_j$ ). So, if  $a$  is a root of  $\text{SRemS}_j$ ,  $j \neq 0$ ,  $\text{SRemS}_{j-1}(a)\text{SRemS}_{j+1}(a) < 0$ , since

$$\text{SRemS}_{j+1} = -\text{SRemS}_{j-1} + \text{Quo}(\text{SRemS}_j, \text{SRemS}_{j-1})\text{SRemS}_j$$

(see Remark 1.4) so that

$$\begin{aligned} & \text{Var}(\text{SRemS}_{j-1}, \text{SRemS}_j, \text{SRemS}_{j+1}; a) \\ &= \text{Var}(\text{SRemS}_{j-1}, \text{SRemS}_j, \text{SRemS}_{j+1}; a') \\ &= 1. \end{aligned}$$

This implies  $\text{Var}(\text{SRemS}(P, Q); a) = \text{Var}(\text{SRemS}(P, Q); a')$ , and similarly  $\text{Var}(\text{SRemS}(P, Q); b) = \text{Var}(\text{SRemS}(P, Q); b')$ .

The proof of the theorem now proceeds by induction on the number  $n \geq 2$  of elements in the signed remainder sequence. The base case  $n=2$  corresponds to  $R=0$  and follows from Lemma 2.59 and Lemma 2.60. Let us suppose that the Theorem holds for  $n-1$  and consider  $P$  and  $Q$  such that their signed remainder sequence has  $n$  elements. The signed remainder sequence of  $Q$  and  $-R$  has  $n-1$  elements and, by the induction hypothesis,

$$\text{Var}(\text{SRemS}(Q, -R); a, b) = \text{Ind}(-R/Q; a, b).$$

So, by Lemma 2.59 and Lemma 2.60,

$$\text{Var}(\text{SRemS}(P, Q); a, b) = \text{Ind}(QP; a, b). \quad \square$$

As a consequence of the above we derive the following theorem.

**Theorem 2.61. [Tarski’s theorem]** *If  $a < b$  are elements of  $\mathbb{R} \cup \{-\infty, +\infty\}$  that are not roots of  $P$ , with  $P, Q \in \mathbb{R}[X]$ , then*

$$\text{Var}(\text{SRemS}(P, P'Q); a, b) = \text{TaQ}(Q, P; a, b).$$

**Proof:** This is immediate from Theorem 2.58 and Proposition 2.57. □

Theorem 2.50 (Sturm’s theorem) is a particular case of Theorem 2.61, taking  $Q = 1$ .

**Proof of Theorem 2.50:** The proof is immediate by taking  $Q = 1$  in Theorem 2.61. □

### 2.3 Projection Theorem for Algebraic Sets

Let  $\mathbb{R}$  be a real closed field. If  $\mathcal{P}$  is a finite subset of  $\mathbb{R}[X_1, \dots, X_k]$ , we write the **set of zeros** of  $\mathcal{P}$  in  $\mathbb{R}^k$  as

$$\text{Zer}(\mathcal{P}, \mathbb{R}^k) = \{x \in \mathbb{R}^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\}.$$

These are the **algebraic sets** of  $\mathbb{R}^k = \text{Zer}(\{0\}, \mathbb{R}^k)$ .

An important way in which this differs from the algebraically closed case is that the common zeros of  $\mathcal{P}$  are also the zeros of a single polynomial  $Q = \sum_{P \in \mathcal{P}} P^2$ .

The smallest family of sets of  $\mathbb{R}^k$  that contains the algebraic sets and is closed under the boolean operations (complementation, finite unions, and finite intersections) is the **constructible sets**.

We define the **semi-algebraic sets** of  $\mathbb{R}^k$  as the smallest family of sets in  $\mathbb{R}^k$  that contains the algebraic sets as well as sets defined by polynomial **inequalities** i.e. sets of the form  $\{x \in \mathbb{R}^k \mid P(x) > 0\}$  for some polynomial  $P \in \mathbb{R}[X_1, \dots, X_k]$ , and which is also closed under the boolean operations (complementation, finite unions, and finite intersections). If the coefficients of the polynomials defining  $S$  lie in a subring  $D \subset \mathbb{R}$ , we say that the semi-algebraic set  $S$  is **defined over  $D$** .

It is obvious that any semi-algebraic set in  $\mathbb{R}^k$  is the finite union of sets of the form  $\{x \in \mathbb{R}^k \mid P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(x) > 0\}$ . These are the **basic semi-algebraic sets**.

Notice that the constructible sets are semi-algebraic as the basic constructible set

$$S = \{x \in \mathbb{R}^k \mid P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(x) \neq 0\}$$



is the basic semi-algebraic set

$$\{x \in \mathbb{R}^k \mid P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q^2(x) > 0\}.$$

The goal of the next pages is to show that the projection of an algebraic set in  $\mathbb{R}^{k+1}$  is a semi-algebraic set of  $\mathbb{R}^k$  if  $\mathbb{R}$  is a real closed field.

This is a new example of the paradigm described in Chapter 1 for extending an algorithm from the univariate case to the multivariate case by viewing the univariate case parametrically. The algebraic set  $Z \subset \mathbb{R}^{k+1}$  can be described as

$$Z = \{(y, x) \in \mathbb{R}^{k+1} \mid P(y, x) = 0\}$$

with  $P \in \mathbb{R}[X_1, \dots, X_k, X_{k+1}]$ , and its projection  $\pi(Z)$  (forgetting the last coordinate) is

$$\pi(Z) = \{y \in \mathbb{R}^k \mid \exists x \in \mathbb{R} P(y, x) = 0\}.$$

For a particular  $y \in \mathbb{R}^k$  we can decide, using Theorem 2.50 (Sturm's theorem) and Remark 2.51, whether or not  $\exists x \in \mathbb{R} P_y(x) = 0$  is true.

Defining  $Z_y = \{x \in \mathbb{R} \mid P_y(x) = 0\}$ , (see Notation 1.18) what is crucial here is to partition the parameter space  $\mathbb{R}^k$  into finitely many parts so that each part is either contained in  $\{y \in \mathbb{R}^k \mid Z_y = \emptyset\}$  or in  $\{y \in \mathbb{R}^k \mid Z_y \neq \emptyset\}$ . Moreover, the algorithm used for constructing the partition ensures that the decision algorithm testing whether  $Z_y$  is empty or not is the same (is uniform) for all  $y$  in any given part. Because of this uniformity, it turns out that each part of the partition is a semi-algebraic set. Since  $\pi(Z)$  is the union of those parts where  $Z_y \neq \emptyset$ ,  $\pi(Z)$  is semi-algebraic being the union of finitely many semi-algebraic sets.

We first introduce some terminology from logic which is useful for the study of semi-algebraic sets.

We define the language of ordered fields by describing the formulas of this language. The definitions are similar to the corresponding notions in Chapter 1, the only difference is the use of inequalities in the atoms. The formulas are built starting with atoms, which are polynomial equations and inequalities. A formula is written using atoms together with the logical connectives "and", "or", and "negation" ( $\wedge$ ,  $\vee$ , and  $\neg$ ) and the existential and universal quantifiers ( $\exists$ ,  $\forall$ ). A formula has free variables, i.e. non-quantified variables, and bound variables, i.e. quantified variables. More precisely, let  $D$  be a subring of  $\mathbb{R}$ . We define the **language of ordered fields with coefficients in  $D$**  as follows. An **atom** is  $P = 0$  or  $P > 0$ , where  $P$  is a polynomial in  $D[X_1, \dots, X_k]$ . We define simultaneously the **formulas** and the set  $\text{Free}(\Phi)$  of **free variables of a formula  $\Phi$**  as follows

- an atom  $P = 0$  or  $P > 0$ , where  $P$  is a polynomial in  $D[X_1, \dots, X_k]$  is a formula with free variables  $\{X_1, \dots, X_k\}$ ,

- if  $\Phi_1$  and  $\Phi_2$  are formulas, then  $\Phi_1 \wedge \Phi_2$  and  $\Phi_1 \vee \Phi_2$  are formulas with  $\text{Free}(\Phi_1 \wedge \Phi_2) = \text{Free}(\Phi_1 \vee \Phi_2) = \text{Free}(\Phi_1) \cup \text{Free}(\Phi_2)$ ,
- if  $\Phi$  is a formula, then  $\neg(\Phi)$  is a formula with  $\text{Free}(\neg(\Phi)) = \text{Free}(\Phi)$ ,
- if  $\Phi$  is a formula and  $X \in \text{Free}(\Phi)$ , then  $(\exists X) \Phi$  and  $(\forall X) \Phi$  are formulas with  $\text{Free}((\exists X) \Phi) = \text{Free}((\forall X) \Phi) = \text{Free}(\Phi) \setminus \{X\}$ .

If  $\Phi$  and  $\Psi$  are formulas,  $\Phi \Rightarrow \Psi$  is the formula  $\neg(\Phi) \vee \Psi$ .

A **quantifier free formula** is a formula in which no quantifier appears, neither  $\exists$  nor  $\forall$ . A **basic formula** is a conjunction of atoms.

The **R-realization of a formula**  $\Phi$  with free variables contained in  $\{Y_1, \dots, Y_k\}$ , denoted  $\text{Reali}(\Phi, \mathbb{R}^k)$ , is the set of  $y \in \mathbb{R}^k$  such that  $\Phi(y)$  is true. It is defined by induction on the construction of the formula, starting from atoms:

$$\begin{aligned} \text{Reali}(P = 0, \mathbb{R}^k) &= \{y \in \mathbb{R}^k \mid P(y) = 0\}, \\ \text{Reali}(P > 0, \mathbb{R}^k) &= \{y \in \mathbb{R}^k \mid P(y) > 0\}, \\ \text{Reali}(P < 0, \mathbb{R}^k) &= \{y \in \mathbb{R}^k \mid P(y) < 0\}, \\ \text{Reali}(\Phi_1 \wedge \Phi_2, \mathbb{R}^k) &= \text{Reali}(\Phi_1, \mathbb{R}^k) \cap \text{Reali}(\Phi_2, \mathbb{R}^k), \\ \text{Reali}(\Phi_1 \vee \Phi_2, \mathbb{R}^k) &= \text{Reali}(\Phi_1, \mathbb{R}^k) \cup \text{Reali}(\Phi_2, \mathbb{R}^k), \\ \text{Reali}(\neg\Phi, \mathbb{R}^k) &= \mathbb{R}^k \setminus \text{Reali}(\Phi, \mathbb{R}^k), \\ \text{Reali}((\exists X) \Phi, \mathbb{R}^k) &= \{y \in \mathbb{R}^k \mid \exists x \in \mathbb{R} \quad (x, y) \in \text{Reali}(\Phi, \mathbb{R}^{k+1})\}, \\ \text{Reali}((\forall X) \Phi, \mathbb{R}^k) &= \{y \in \mathbb{R}^k \mid \forall x \in \mathbb{R} \quad (x, y) \in \text{Reali}(\Phi, \mathbb{R}^{k+1})\} \end{aligned}$$

Two formulas  $\Phi$  and  $\Psi$  such that  $\text{Free}(\Phi) = \text{Free}(\Psi) = \{Y_1, \dots, Y_k\}$  are **R-equivalent** if  $\text{Reali}(\Phi, \mathbb{R}^k) = \text{Reali}(\Psi, \mathbb{R}^k)$ . If there is no ambiguity, we simply write  $\text{Reali}(\Phi)$  for  $\text{Reali}(\Phi, \mathbb{R}^k)$  and talk about realization and equivalence.

It is clear that a set is semi-algebraic if and only if it can be represented as the realization of a quantifier free formula. It is also easy to see that any formula in the language of fields with coefficients in  $\mathbb{D}$  is R-equivalent to a formula

$$\Phi(Y) = (\text{Qu}_1 X_1) \dots (\text{Qu}_m X_m) \mathcal{B}(X_1, \dots, X_m, Y_1, \dots, Y_k)$$

where each  $\text{Qu}_i \in \{\forall, \exists\}$  and  $\mathcal{B}$  is a quantifier free formula involving polynomials in  $\mathbb{D}[X_1, \dots, X_m, Y_1, \dots, Y_k]$ . This is called its **prenex normal form** (see Section 10, Chapter 1 of [115]). The variables  $X_1, \dots, X_m$  are called **bound variables**. If a formula has no free variables, then it is called a **sentence**, and is either R-equivalent to true, when  $\text{Reali}(\Phi, \{0\}) = \{0\}$ , or R-equivalent to false, when  $\text{Reali}(\Phi, \{0\}) = \emptyset$ . For example,  $1 > 0$  is R-equivalent to true, and  $1 < 0$  is R-equivalent to false.

We now prove that the projection of an algebraic set is semi-algebraic.

**Theorem 2.62.** *Given an algebraic set of  $\mathbb{R}^{k+1}$  defined over  $D$ , its projection to  $\mathbb{R}^k$  is a semi-algebraic set defined over  $D$ .*

Before proving Theorem 2.62, let us explain the mechanism of its proof on an example.

*Example 2.63.* We describe the projection of the algebraic set

$$\{(a, b, c, X) \in \mathbb{R}^4 \mid X^4 + aX^2 + bX + c = 0\}$$

to  $\mathbb{R}^3$ , i.e. the set

$$\{(a; b; c) \in \mathbb{R}^3 \mid \exists X \in \mathbb{R} \quad X^4 + aX^2 + bX + c = 0\},$$

as a semi-algebraic set.

We look at all leaves of  $\text{TRems}(P, P')$  and at all possible signs for leading coefficients of all possible signed pseudo-remainders (using Example 1.15). We denote by  $n$  the difference between the number of sign variations at  $-\infty$  and  $+\infty$  in the Sturm sequence of  $P = X^4 + aX^2 + bX + c$  for each case. We indicate for each leaf  $L$  of  $\text{TRems}(P, P')$  the basic formula  $\mathcal{C}_L$  and the degrees occurring in the signed pseudo-remainder sequence of  $P$  and  $P'$  along the path  $\mathcal{B}_L$ .

$$(a \neq 0 \wedge s \neq 0 \wedge \delta \neq 0, (4, 3, 2, 1, 0))$$

$a$	-	-	-	-	+	+	+	+
$s$	+	+	-	-	+	+	-	-
$\delta$	+	-	+	-	+	-	+	-
$n$	4	2	0	2	0	-2	0	2

The first column can be read as follows: for every polynomial

$$P = X^4 + aX^2 + bX + c$$

satisfying  $a < 0$ ,  $s > 0$ ,  $\delta > 0$ , the number of real roots is 4. Indeed the leading coefficients of the signed pseudo-remainder sequence of  $P$  and  $P'$  are  $1, 4, -a, 64s, 16384a^2\delta$  (see Example 1.17) and the degrees of the polynomials in the signed pseudo-remainder sequence of  $P$  and  $P'$  are  $4, 3, 2, 1, 0$ , the signs of the signed pseudo-remainder sequence of  $P$  and  $P'$  at  $-\infty$  are  $+ - + - +$  and at  $+\infty$  are  $+ + + + +$ . We can apply Theorem 2.50 (Sturm's Theorem).

The other columns can be read similarly. Notice that  $n$  can be negative (for  $a > 0, s > 0, \delta < 0$ ). Though this looks paradoxical, Sturm's theorem is not violated. It only means that there is no polynomial  $P \in \mathbb{R}[X]$  with  $P = X^4 + aX^2 + bX + c$  and  $a > 0, s > 0, \delta < 0$ . Notice that even when  $n$  is non-negative, there might be no polynomial  $P \in \mathbb{R}[X]$  with  $P = X^4 + aX^2 + bX + c$  and  $(a, s, \delta)$  satisfying the corresponding sign condition.

Similarly, for the other leaves of  $\text{TRems}(P, P')$

$$(a \neq 0 \wedge s \neq 0 \wedge \delta = 0, (4, 3, 2, 1))$$

$$\frac{\begin{array}{c|cccc} a & - & - & + & + \\ s & + & - & + & - \\ \hline n & 3 & 1 & -1 & 1 \end{array}}$$

$$(a \neq 0 \wedge s = 0 \wedge t \neq 0, (4, 3, 2, 0))$$

$$\frac{\begin{array}{c|cccc} a & - & - & + & + \\ t & + & - & + & - \\ \hline n & 2 & 2 & 0 & 0 \end{array}}$$

$$(a \neq 0 \wedge s = t = 0, (4, 3, 2))$$

$$\frac{\begin{array}{c|cc} a & - & + \\ \hline n & 2 & 0 \end{array}}$$

$$(a = 0 \wedge b \neq 0 \wedge u \neq 0, (4, 3, 1, 0))$$

$$\frac{\begin{array}{c|cccc} b & + & + & - & - \\ u & + & - & + & - \\ \hline n & 2 & 0 & 0 & 2 \end{array}}$$

$$(a = 0 \wedge b \neq 0 \wedge u = 0, (4, 3, 1))$$

$$\frac{\begin{array}{c|cc} b & + & - \\ \hline n & 1 & 1 \end{array}}$$

$$(a = b = 0 \wedge c \neq 0, (4, 3, 0))$$

$$\frac{\begin{array}{c|cc} c & + & - \\ \hline n & 0 & 2 \end{array}}$$

$$(a = b = c = 0, (4, 3))$$

$$n = 1$$

Finally, the formula  $\exists X \quad X^4 + a X^2 + b X + c = 0$  is R-equivalent to the quantifier-free formula  $\Phi(a, b, c)$ :

$$\begin{aligned}
& (a < 0 \wedge s > 0) \\
& \vee (a < 0 \wedge s < 0 \wedge \delta < 0) \\
& \vee (a > 0 \wedge s < 0 \wedge \delta < 0) \\
& \vee (a < 0 \wedge s \neq 0 \wedge \delta = 0) \\
& \vee (a > 0 \wedge s < 0 \wedge \delta = 0) \\
& \vee (a < 0 \wedge s = 0 \wedge t \neq 0) \\
& \vee (a < 0 \wedge s = 0 \wedge t = 0) \\
& \vee (a = 0 \wedge b < 0 \wedge u < 0) \\
& \vee (a = 0 \wedge b > 0 \wedge u > 0) \\
& \vee (a = 0 \wedge b \neq 0 \wedge u = 0) \\
& \vee (a = 0 \wedge b = 0 \wedge c < 0) \\
& \vee (a = 0 \wedge b = 0 \wedge c = 0),
\end{aligned}$$

by collecting all the sign conditions with  $n \geq 1$ . Thus, we have proven that the projection of the algebraic set

$$\{(x, a, b, c) \in \mathbb{R}^4 \mid x^4 + a x^2 + b x + c\}$$

into  $\mathbb{R}^3$  is the semi-algebraic subset  $\text{Reali}(\Phi, \mathbb{R}^3)$ .  $\square$

The proof of Theorem 2.62 follows closely the method illustrated in the example.

**Proof of Theorem 2.62:** Let  $Z = \{z \in \mathbb{R}^{k+1} \mid P(z) = 0\}$ . Let  $Z'$  be the intersection of  $Z$  with the subset of  $(y, x) \in \mathbb{R}^{k+1}$  such that  $P_y$  is not identically zero.

Let  $L$  be a leaf of  $\text{TRems}(P, P')$ , and let  $\mathcal{A}(L)$  be the set of non-zero polynomials in  $D[Y_1, \dots, Y_k]$  appearing in the basic formula  $\mathcal{C}_L$ , (see Notation 1.18).

Let  $\mathcal{L}$  be the set of all leaves of  $\text{TRems}(P, P')$ , and

$$\mathcal{A} = \bigcup_{L \in \mathcal{L}} \mathcal{A}(L) \subset D[Y_1, \dots, Y_k].$$

If  $\tau \in \{0, 1, -1\}^{\mathcal{A}}$ , we define the realization of  $\tau$  by

$$\text{Reali}(\tau) = \{y \in \mathbb{R}^k \mid \bigwedge_{A \in \mathcal{A}} \text{sign}(A(y)) = \tau(A)\}.$$

Let  $Z_y = \{x \in \mathbb{R} \mid P(y, x) = 0\}$ . Note that  $\text{Reali}(\tau) \subset \{y \in \mathbb{R}^k \mid Z_y \neq \emptyset\}$  or  $\text{Reali}(\tau) \subset \{y \in \mathbb{R}^k \mid Z_y = \emptyset\}$ , by Theorem 2.50 (Sturm's theorem) and Remark 2.51. Let

$$\Sigma = \{\tau \in \{0, 1, -1\}^{\mathcal{A}} \mid \forall y \in \text{Reali}(\tau) \quad Z_y \neq \emptyset\}.$$

It is clear that the semi-algebraic set  $\bigcup_{\tau \in \Sigma} \text{Reali}(\tau)$  coincides with the projection of  $S'$ .

The fact that the projection of the intersection of  $Z$  with the subset of  $(y, x) \in \mathbb{R}^{k+1}$  such that  $P_y$  is identically zero is semi-algebraic is obvious.

Thus the whole projection of  $Z = Z' \cup (Z \setminus Z')$  is semi-algebraic since it is a union of semi-algebraic sets. □

## 2.4 Projection Theorem for Semi-Algebraic Sets

We are going to prove by a similar method that the projection of a semi-algebraic set is semi-algebraic. We start with a decision algorithm deciding if a given sign condition has a non-empty realization at the zeroes of a univariate polynomial.

When  $P$  and  $Q$  have no common roots, we can find the number of roots of  $P$  at each possible sign of  $Q$  in terms of the Tarski-queries of 1 and  $Q$  for  $P$ .

We denote

$$\begin{aligned} Z &= \text{Zer}(P, \mathbb{R}) \\ &= \{x \in \mathbb{R} \mid P(x) = 0\}, \\ \text{Reali}(Q = 0, Z) &= \{x \in Z \mid \text{sign}(Q(x)) = 0\} = \{x \in Z \mid Q(x) = 0\}, \\ \text{Reali}(Q > 0, Z) &= \{x \in Z \mid \text{sign}(Q(x)) = 1\} = \{x \in Z \mid Q(x) > 0\}, \\ \text{Reali}(Q < 0, Z) &= \{x \in Z \mid \text{sign}(Q(x)) = -1\} = \{x \in Z \mid Q(x) < 0\}, \end{aligned}$$

and  $c(Q = 0, Z)$ ,  $c(Q > 0, Z)$ ,  $c(Q < 0, Z)$  are the cardinalities of the corresponding sets.

**Proposition 2.64.** *If  $P$  and  $Q$  have no common roots in  $\mathbb{R}$ , then*

$$\begin{aligned} c(Q > 0, Z) &= (\text{TaQ}(1, P) + \text{TaQ}(Q, P))/2, \\ c(Q < 0, Z) &= (\text{TaQ}(1, P) - \text{TaQ}(Q, P))/2. \end{aligned}$$

**Proof:** We have

$$\begin{aligned} \text{TaQ}(1, P) &= c(Q > 0, Z) + c(Q < 0, Z), \\ \text{TaQ}(Q, P) &= c(Q > 0, Z) - c(Q < 0, Z). \end{aligned}$$

Now solve. □

With a little more effort, we can find the number of roots of  $P$  at each possible sign of  $Q$  in terms of the Tarski-queries of 1,  $Q$ , and  $Q^2$  for  $P$ .

**Proposition 2.65.** *The following holds*

$$\begin{aligned} c(Q=0, Z) &= \text{TaQ}(1, P) - \text{TaQ}(Q^2, P), \\ c(Q>0, Z) &= (\text{TaQ}(Q^2, P) + \text{TaQ}(Q, P))/2, \\ c(Q<0, Z) &= (\text{TaQ}(Q^2, P) - \text{TaQ}(Q, P))/2. \end{aligned}$$

**Proof:** Indeed, we have

$$\begin{aligned} \text{TaQ}(1, P) &= c(Q=0, Z) + c(Q>0, Z) + c(Q<0, Z), \\ \text{TaQ}(Q, P) &= c(Q>0, Z) - c(Q<0, Z), \\ \text{TaQ}(Q^2, P) &= c(Q>0, Z) + c(Q<0, Z). \end{aligned}$$

Now solve. □

We want to extend these results to the case of many polynomials.

We consider a  $P \in \mathbb{R}[X]$  with  $P$  not identically zero,  $\mathcal{Q}$  a finite subset of  $\mathbb{R}[X]$ , and the finite set  $Z = \text{Zer}(P, \mathbb{R}) = \{x \in \mathbb{R} \mid P(x) = 0\}$ .

We will give an expression for the number of elements of  $Z$  at which  $\mathcal{Q}$  satisfies a given sign condition  $\sigma$ .

Let  $\sigma$  be a sign condition on  $\mathcal{Q}$  i.e. an element of  $\{0, 1, -1\}^{\mathcal{Q}}$ . The **realization of the sign condition  $\sigma$  over  $Z$**  is

$$\text{Reali}(\sigma, Z) = \{x \in \mathbb{R} \mid P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(x)) = \sigma(Q)\}.$$

Its cardinality is denoted  $c(\sigma, Z)$ .

Given  $\alpha \in \{0, 1, 2\}^{\mathcal{Q}}$ , and  $\sigma \in \{0, 1, -1\}^{\mathcal{Q}}$  we write  $\sigma^\alpha$  for  $\prod_{Q \in \mathcal{Q}} \sigma(Q)^{\alpha(Q)}$ , and  $\mathcal{Q}^\alpha$  for  $\prod_{Q \in \mathcal{Q}} Q^{\alpha(Q)}$ . When  $\text{Reali}(\sigma, Z) \neq \emptyset$ , the sign of  $\mathcal{Q}^\alpha$  is fixed on  $\text{Reali}(\sigma, Z)$  and is equal to  $\sigma^\alpha$ , with the convention that  $0^0 = 1$ .

We number the elements of  $\mathcal{Q}$  so that  $\mathcal{Q} = \{Q_1, \dots, Q_s\}$  and use the lexicographical orderings on  $\{0, 1, 2\}^{\mathcal{Q}}$  (with  $0 < 1 < 2$ ) and  $\{0, 1, -1\}^{\mathcal{Q}}$  (with  $0 < 1 < -1$ ) (see Definition 2.14).

Given a list of elements  $A = \alpha_1, \dots, \alpha_m$  of  $\{0, 1, 2\}^{\mathcal{Q}}$  with  $\alpha_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_m$ , we define

$$\begin{aligned} \mathcal{Q}^A &= \mathcal{Q}^{\alpha_1}, \dots, \mathcal{Q}^{\alpha_m} \\ \text{TaQ}(\mathcal{Q}^A, P) &= \text{TaQ}(\mathcal{Q}^{\alpha_1}, P), \dots, \text{TaQ}(\mathcal{Q}^{\alpha_m}, P). \end{aligned}$$

Given a list of elements  $\Sigma = \sigma_1, \dots, \sigma_n$  of  $\{0, 1, -1\}^{\mathcal{Q}}$ , with  $\sigma_1 <_{\text{lex}} \dots <_{\text{lex}} \sigma_n$ , we define

$$\begin{aligned} \text{Reali}(\Sigma, Z) &= \text{Reali}(\sigma_1, Z), \dots, \text{Reali}(\sigma_n, Z) \\ c(\Sigma, Z) &= c(\sigma_1, Z), \dots, c(\sigma_n, Z). \end{aligned}$$

**Definition 2.66.** The **matrix of signs of  $\mathcal{Q}^A$  on  $\Sigma$**  is the  $m \times n$  matrix  $\text{Mat}(A, \Sigma)$  whose  $i, j$ -th entry is  $\sigma_j^{\alpha_i}$ . □

*Example 2.67.* If  $\mathcal{Q} = \{Q_1, Q_2\}$  and  $A = \{0, 1, 2\}^{\{Q_1, Q_2\}}$ ,  $\{Q_1, Q_2\}^A$  is the list  $1, Q_2, Q_2^2, Q_1, Q_1Q_2, Q_1Q_2^2, Q_1^2, Q_1^2Q_2, Q_1^2Q_2^2$ . Taking  $\Sigma = \{0, 1, -1\}^{\{Q_1, Q_2\}}$ , i.e. the list

$$\begin{aligned} &Q_1 = 0 \wedge Q_2 = 0, Q_1 = 0 \wedge Q_2 > 0, Q_1 = 0 \wedge Q_2 < 0, \\ &Q_1 > 0 \wedge Q_2 = 0, Q_1 > 0 \wedge Q_2 > 0, Q_1 > 0 \wedge Q_2 < 0, \\ &Q_1 < 0 \wedge Q_2 = 0, Q_1 < 0 \wedge Q_2 > 0, Q_1 < 0 \wedge Q_2 < 0, \end{aligned}$$

the matrix of signs of these nine polynomials on these nine sign conditions is

$$\text{Mat}(A, \Sigma) = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 \\ 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & -1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & -1 & -1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 1 \end{bmatrix}.$$

For example, the 5-th row of the matrix reads as follows: the signs of the 5-th polynomial of  $\mathcal{Q}^A$  which is  $Q_1Q_2$  on the 9 sign conditions of  $\Sigma$  are

$$[ 0 \ 0 \ 0 \ 0 \ 1 \ -1 \ 0 \ -1 \ 1 ]. \quad \square$$

**Proposition 2.68.** *If  $\bigcup_{\sigma \in \Sigma} \text{Reali}(\sigma, Z) = Z$  then*

$$\text{Mat}(A, \Sigma) \cdot c(\Sigma, Z) = \text{TaQ}(\mathcal{Q}^A, P).$$

**Proof:** It is obvious since the  $(i, j)$ -th entry of  $\text{Mat}(A, \Sigma)$  is  $\sigma_j^{\alpha_i}$ . □

Note that when  $\mathcal{Q} = \{Q\}$ ,  $A = \{0, 1, 2\}^{\{Q\}}$  and  $\Sigma = \{0, 1, -1\}^{\{Q\}}$  the conclusion of Proposition 2.68 is

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(Q=0, Z) \\ c(Q>0, Z) \\ c(Q<0, Z) \end{bmatrix} = \begin{bmatrix} \text{TaQ}(1, P) \\ \text{TaQ}(Q, P) \\ \text{TaQ}(Q^2, P) \end{bmatrix} \quad (2.6)$$

which was hidden in the proof of Proposition 2.65.

It follows from Proposition 2.68 that when the matrix  $M(\mathcal{Q}^A, \Sigma)$  is invertible, we can express  $c(\Sigma, Z)$  in terms of  $\text{TaQ}(\mathcal{Q}^A, P)$ . This is the case when  $A = \{0, 1, 2\}^{\mathcal{Q}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{Q}}$ , as we will see now.

**Notation 2.69. [Tensor product]** Let  $M$  and  $M' = [m'_{ij}]$  be two matrices with respective dimensions  $n \times m$  and  $n' \times m'$ . The matrix  $M \otimes M'$  is the  $nn' \times mm'$  matrix

$$[ m_{ij}m'_{ij} ].$$

The matrix  $M \otimes M'$  is the **tensor product** of  $M$  and  $M'$ . □



*Example 2.70.* If

$$M = M' = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix},$$

$$M \otimes M' = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 \\ 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & -1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & -1 & -1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 1 \end{bmatrix}.$$

Notice that  $M \otimes M'$  coincides with the matrix of signs of  $A = \{0, 1, 2\}^{\{Q_1, Q_2\}}$  on  $\Sigma = \{0, 1, -1\}^{\{Q_1, Q_2\}}$ .  $\square$

**Notation 2.71.** Let  $M_s$  be the  $3^s \times 3^s$  matrix defined inductively by

$$M_1 = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix}$$

$$M_{t+1} = M_t \otimes M_1.$$

$\square$

**Exercise 2.13.** Prove that  $M_s$  is invertible using induction on  $s$ .

**Proposition 2.72.** Let  $\mathcal{Q}$  be a finite set of polynomials with  $s$  elements,  $A = \{0, 1, 2\}^{\mathcal{Q}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{Q}}$ , ordered lexicographically. Then

$$\text{Mat}(A, \Sigma) = M_s.$$

**Proof:** The proof is by induction on  $s$ . If  $s=1$ , the claim is Equation (2.6). If the claim holds for  $s$ , it holds also for  $s+1$  given the definitions of  $M_{s+1}$ , of  $\text{Mat}(A, \Sigma)$ , and the orderings on  $A = \{0, 1, 2\}^{\mathcal{Q}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{Q}}$ .  $\square$

So, Proposition 2.68 and Proposition 2.72 imply

**Corollary 2.73.**  $M_s \cdot c(\Sigma, Z) = \text{TaQ}(\mathcal{Q}^A, P)$ .

We now have all the ingredients needed to decide whether a subset of  $\mathbb{R}$  defined by a sign condition is empty or not, with the following two lemmas.

**Lemma 2.74.** Let  $Z = \text{Zer}(P, \mathbb{R})$  be a finite set and let  $\sigma$  be a sign condition on  $\mathcal{Q}$ . Whether or not  $\text{Reali}(\sigma, Z) = \emptyset$  is determined by the degrees of the polynomials in the signed pseudo-remainder sequences of  $P$ ,  $P'Q^\alpha$  and the signs of their leading coefficients for all  $\alpha \in A = \{0, 1, 2\}^{\mathcal{Q}}$ .

**Proof:** For each  $\alpha \in \{0, 1, 2\}^{\mathcal{Q}}$ , the degrees and the signs of the leading coefficients of all of the polynomials in the signed pseudo-remainder sequences  $\text{SRemS}(P, P'Q^\alpha)$  clearly determine the number of sign variations of  $\text{SRemS}(P, P'Q^\alpha)$  at  $-\infty$  and  $+\infty$ , i.e.  $\text{Var}(\text{SRemS}(P, P'Q^\alpha); -\infty)$  and  $\text{Var}(\text{SRemS}(P, P'Q^\alpha); +\infty)$ , and their difference is  $\text{TaQ}(Q^\alpha, P)$  by Theorem 2.61. Using Propositions 2.72, Proposition 2.68, and Corollary 2.73

$$M_s^{-1} \cdot \text{TaQ}(Q^A, P) = c(\Sigma, Z).$$

Denoting the row of  $M_s^{-1}$  that corresponds to the row of  $\sigma$  in  $c(\Sigma, Z)$  by  $r_\sigma$ , we see that  $r_\sigma \cdot \text{TaQ}(Q^A, P) = c(\sigma, Z)$ . Finally,

$$\text{Reali}(\sigma, Z) = \{x \in \mathbb{R} \mid P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(x)) = \sigma(Q)\}$$

is non-empty if and only if  $c(\sigma, Z) > 0$ . □

**Lemma 2.75.** *Let  $\sigma$  be a strict sign condition on  $\mathcal{Q}$ . Whether or not  $\text{Reali}(\sigma) = \emptyset$  is determined by the degrees and the signs of the leading coefficients of the polynomials in  $\text{Var}(\text{SRemS}(C, C'))$  (with  $C = \prod_{Q \in \mathcal{Q}} Q$ ) and the signs of the leading coefficients of the polynomials in  $\text{Var}(\text{SRemS}(C', C''Q^\alpha))$  for all  $\alpha \in A = \{0, 1, 2\}^{\mathcal{Q}}$ .*

**Proof:** Recall (Theorem 2.50) that the number of roots of  $C$  is determined by the signs of the leading coefficients of  $\text{Var}(\text{SRemS}(C, C'))$ .

- If  $C$  has no roots, then each  $Q \in \mathcal{Q}$  has constant sign which is the same as the sign of its leading coefficient.
- If  $C$  has one root, then the possible sign conditions on  $\mathcal{Q}$  are determined by the sign conditions on  $\mathcal{Q}$  at  $+\infty$  and at  $-\infty$ .
- If  $C$  has at least two roots, then all intervals between two roots of  $C$  contain a root of  $C'$  and thus all sign conditions on  $\mathcal{Q}$  are determined by the sign conditions on  $\mathcal{Q}$  at  $+\infty$  and at  $-\infty$  and by the sign conditions on  $\mathcal{Q}$  at the roots of  $C'$ . This is covered by Lemma 2.74. □

The goal of the remainder of the section is to show that the semi-algebraic sets in  $\mathbb{R}^{k+1}$  are closed under projection if  $\mathbb{R}$  is a real closed field. The result is a generalization of Theorem 2.62 and the proof is based on a similar method.

Let us now describe our algorithm for proving that the projection of a semi-algebraic set is semi-algebraic. Using how to decide whether or not a basic semi-algebraic set in  $\mathbb{R}$  is empty (see Lemmas 2.74 and 2.75), we can show that the projection from  $\mathbb{R}^{k+1}$  to  $\mathbb{R}^k$  of a basic semi-algebraic set is semi-algebraic. This is a new example of our paradigm for extending an algorithm from the univariate case to the multivariate case by viewing the univariate case parametrically. The basic semi-algebraic set  $S \subset \mathbb{R}^{k+1}$  can be described as

$$S = \{x \in \mathbb{R}^{k+1} \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q(x) > 0\}$$

with  $\mathcal{P}, \mathcal{Q}$  finite subsets of  $\mathbb{R}[X_1, \dots, X_k, X_{k+1}]$ , and its projection  $\pi(S)$  (forgetting the last coordinate) is

$$\pi(S) = \{y \in \mathbb{R}^k \mid \exists x \in \mathbb{R} \left( \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \bigwedge_{Q \in \mathcal{Q}} Q_y(x) > 0 \right)\}.$$

For a particular  $y \in \mathbb{R}^k$  we can decide, using Lemmas 2.74 and 2.75, whether or not

$$\exists x \in \mathbb{R} \left( \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \bigwedge_{Q \in \mathcal{Q}} Q_y(x) > 0 \right)$$

is true.

What is crucial here is to partition the parameter space  $\mathbb{R}^k$  into finitely many parts so that each part is either contained in  $\{y \in \mathbb{R}^k \mid S_y = \emptyset\}$  or in  $\{y \in \mathbb{R}^k \mid S_y \neq \emptyset\}$ , where

$$S_y = \{x \in \mathbb{R} \mid \bigwedge_{P \in \mathcal{P}} P_y(x) = 0 \wedge \bigwedge_{Q \in \mathcal{Q}} Q_y(x) > 0\}.$$

Moreover, the algorithm used for constructing the partition ensures that the decision algorithm testing whether  $S_y$  is empty or not is the same (is uniform) for all  $y$  in any given part. Because of this uniformity, it turns out that each part of the partition is a semi-algebraic set. Since  $\pi(S)$  is the union of those parts where  $S_y \neq \emptyset$ ,  $\pi(S)$  is semi-algebraic being the union of finitely many semi-algebraic sets.

**Theorem 2.76. [Projection theorem for semi-algebraic sets]** *Given a semi-algebraic set of  $\mathbb{R}^{k+1}$  defined over  $D$ , its projection to  $\mathbb{R}^k$  is a semi-algebraic set defined over  $D$ .*

**Proof:** Since every semi-algebraic set is a finite union of basic semi-algebraic sets it is sufficient to prove that the projection of a basic semi-algebraic set is semi-algebraic. Suppose that the basic semi-algebraic set  $S$  in  $\mathbb{R}^{k+1}$  is

$$\text{Reali}(\sigma, Z) = \{(y, x) \in \mathbb{R}^k \times \mathbb{R} \mid P(y, x) = \bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(y, x)) = \sigma(Q)\},$$

with  $Z = \{z \in \mathbb{R}^{k+1} \mid P(z) = 0\}$ . Let  $S'$  be the intersection of  $S$  with the subset of  $(y, x) \in \mathbb{R}^{k+1}$  such that  $P_y$  is not identically zero.

Let  $L$  be a function on  $\{0, 1, 2\}^{\mathcal{Q}}$  associating to each  $\alpha \in \{0, 1, 2\}^{\mathcal{Q}}$  a leaf  $L_\alpha$  of  $\text{TRems}(P, P'Q^\alpha)$ , and let  $\mathcal{A}(L_\alpha)$  be the set of non-zero polynomials in  $D[Y_1, \dots, Y_k]$  appearing in the quantifier free formula  $\mathcal{C}_{L_\alpha}$ , (see Notation 1.18).

Let  $\mathcal{L}$  be the set of all functions  $L$  on  $\{0, 1, 2\}^{\mathcal{Q}}$  associating to each  $\alpha$  a leaf  $L_\alpha$  of  $\text{TRems}(P, P'Q^\alpha)$ , and

$$\mathcal{A} = \bigcup_{L \in \mathcal{L}} \bigcup_{\alpha \in \{0,1,2\}^{\mathcal{Q}}} \mathcal{A}(L_\alpha) \subset D[Y_1, \dots, Y_k].$$

Note that since  $\mathcal{A}$  contains the coefficients of  $P'$ , the signs of the coefficients of  $P$  are fixed as soon as the signs of the polynomials in  $\mathcal{A}$  are fixed.

If  $\tau \in \{0, 1, -1\}^{\mathcal{A}}$ , we define the realization of  $\tau$  by

$$\text{Reali}(\tau) = \{y \in \mathbb{R}^k \mid \bigwedge_{A \in \mathcal{A}} \text{sign}(A(y)) = \tau(A)\}.$$

Let  $Z_y = \{x \in \mathbb{R} \mid P(y, x) = 0\}$ ,  $\sigma_y(Q_y) = \sigma(Q)$ , and note that either

$$\text{Reali}(\tau) \subset \{y \in \mathbb{R}^k \mid \text{Reali}(\sigma_y, Z_y) \neq \emptyset\}$$

or

$$\text{Reali}(\tau) \subset \{y \in \mathbb{R}^k \mid \text{Reali}(\sigma_y, Z_y) = \emptyset\},$$

by Lemma 2.74. Let

$$\Sigma = \{\tau \in \{0, 1, -1\}^{\mathcal{A}} \mid \forall y \in \text{Reali}(\tau) \quad \text{Reali}(\sigma_y, Z_y) \neq \emptyset\}.$$

It is clear that the semi-algebraic set  $\bigcup_{\tau \in \Sigma} \text{Reali}(\tau)$  coincides with the projection of  $S'$ .

The fact that the projection of the intersection of  $S$  with the subset of  $(y, x) \in \mathbb{R}^{k+1}$  such that  $P_y$  is identically zero is semi-algebraic follows in a similar way, using Lemma 2.75.

Thus the whole projection  $S = S' \cup (S \setminus S')$  is semi-algebraic as a union of semi-algebraic sets. □

**Exercise 2.14.** Find the conditions on  $a, b$  such that  $X^3 + aX + b$  has a strictly positive real root.

## 2.5 Applications

### 2.5.1 Quantifier Elimination and the Transfer Principle

As in Chapter 1, the projection theorem (Theorem 2.76) implies that the theory of real closed fields admits quantifier elimination in the language of ordered fields, which is the following theorem.

**Theorem 2.77. [Quantifier Elimination over Real Closed Fields]**

*Let  $\Phi(Y)$  be a formula in the language of ordered fields with coefficients in an ordered ring  $D$  contained in the real closed field  $\mathbb{R}$ . Then there is a quantifier free formula  $\Psi(Y)$  with coefficients in  $D$  such that for every  $y \in \mathbb{R}^k$ , the formula  $\Phi(y)$  is true if and only if the formula  $\Psi(y)$  is true.*

The proof of the theorem is by induction on the number of quantifiers, using as base case the elimination of an existential quantifier which is given by Theorem 2.76.

**Proof:** Given a formula  $\Theta(Y) = (\exists X)\mathcal{B}(X, Y)$ , where  $\mathcal{B}$  is a quantifier free formula whose atoms are equations and inequalities involving polynomials in  $D[X, Y_1, \dots, Y_k]$ , Theorem 2.76 shows that there is a quantifier free formula  $\Xi(Y)$  whose atoms are equations and inequalities involving polynomials in  $D[X, Y_1, \dots, Y_k]$  and that is equivalent to  $\Theta(Y)$ . This is because  $\text{Reali}(\Theta(Y), \mathbb{R}^k)$  which is the projection of the semi-algebraic set  $\text{Reali}(\mathcal{B}(X, Y), \mathbb{R}^{k+1})$  defined over  $D$  is semi-algebraic and defined over  $D$ , and semi-algebraic sets defined over  $D$  are realizations of quantifier free formulas with coefficients in  $D$ . Since  $(\forall X)\Phi$  is equivalent to  $\neg((\exists X)\neg(\Phi))$ , the theorem immediately follows by induction on the number of quantifiers.  $\square$

**Corollary 2.78.** *Let  $\Phi(Y)$  be a formula in the language of ordered fields with coefficients in  $D$ . The set  $\{y \in \mathbb{R}^k \mid \Phi(y)\}$  is semi-algebraic.*

**Corollary 2.79.** *A subset of  $\mathbb{R}$  defined by a formula in the language of ordered fields with coefficients in  $\mathbb{R}$  is a finite union of points and intervals.*

**Proof:** By Theorem 2.77 a subset of  $\mathbb{R}$  defined by a formula in the language of ordered fields with coefficients in  $\mathbb{R}$  is semi-algebraic and this is clearly a finite union of points and intervals.  $\square$

**Exercise 2.15.** Show that the set  $\{(x, y) \in \mathbb{R}^2 \mid \exists n \in \mathbb{N} \quad y = nx\}$  is not a semi-algebraic set.

Theorem 2.77 immediately implies the following theorem known as the Tarski-Seidenberg Principle or the Transfer Principle for real closed fields.

**Theorem 2.80. [Tarski-Seidenberg principle]** *Suppose that  $\mathbb{R}'$  is a real closed field that contains the real closed field  $\mathbb{R}$ . If  $\Phi$  is a sentence in the language of ordered fields with coefficients in  $\mathbb{R}$ , then it is true in  $\mathbb{R}$  if and only if it is true in  $\mathbb{R}'$ .*

**Proof:** By Theorem 2.77, there is a quantifier free formula  $\Psi$   $\mathbb{R}$ -equivalent to  $\Phi$ . It follows from the proof of Theorem 2.76 that  $\Psi$  is  $\mathbb{R}'$ -equivalent to  $\Phi$  as well. Notice, too, that  $\Psi$  is a boolean combination of atoms of the form  $c = 0, c > 0$ , or  $c < 0$ , where  $c \in \mathbb{R}$ . Clearly,  $\Psi$  is true in  $\mathbb{R}$  if and only if it is true in  $\mathbb{R}'$ .  $\square$

Since any real closed field contains the real closure of  $\mathbb{Q}$ , a consequence of Theorem 2.80 is

**Theorem 2.81.** *Let  $\mathbb{R}$  be a real closed field. A sentence in the language of fields with coefficients in  $\mathbb{Q}$  is true in  $\mathbb{R}$  if and only if it is true in any real closed field.*

The following application of quantifier elimination will be useful later in the book.

**Proposition 2.82.** *Let  $F$  be an ordered field and  $R$  its real closure. A semi-algebraic set  $S \subset R^k$  can be defined by a quantifier free formula with coefficients in  $F$ .*

**Proof:** Any element  $a \in R$  is algebraic over  $F$ , and is thus a root of a polynomial  $P_a(X) \in F[X]$ . Suppose that  $a = a_j$  where  $a_1 < \dots < a_\ell$  are the roots of  $P_a$  in  $R$ .

Let  $\Delta_a(Y)$  be the formula

$$\begin{aligned} & (\exists Y_1) \dots (\exists Y_\ell) [Y_1 < Y_2 < \dots < Y_\ell \wedge (P_a(Y_1) = \dots = P_a(Y_\ell) = 0) \\ & \wedge ((\forall X) P_a(X) = 0 \Rightarrow (X = Y_1 \vee \dots \vee X = Y_\ell)) \wedge Y = Y_j]. \end{aligned}$$

Then, for  $y \in R$ ,  $\Delta_a(y)$  is true if and only if  $y = a$ .

Let  $A$  be the finite set of elements of  $R \setminus F$  appearing in a quantifier free formula  $\Phi$  with coefficients in  $R$  such that  $S = \{x \in R^k \mid \Phi(x)\}$ . For each  $a \in A$ , replacing each occurrence of  $a$  in  $\Phi$  by new variables  $Y_a$  gives a formula  $\Psi(X, Y)$ , with  $Y = (Y_a, a \in A)$ . Denoting  $n = \#(A)$ , it is clear that  $S = \{x \in R^k \mid \forall y \in R^n (\bigwedge_{a \in A} \Delta_a(y_a) \Rightarrow \Psi(x, y))\}$ .

The conclusions follows from Theorem 2.77 since the formula

$$\forall Y \left( \bigwedge_{a \in A} \Delta_a(Y_a) \Rightarrow \Psi(X, Y) \right)$$

is equivalent to a quantifier free formula with coefficients in  $F$ . □

### 2.5.2 Semi-Algebraic Functions

Since the main objects of our interest are the semi-algebraic sets we want to introduce mappings which preserve semi-algebraicity. These are the semi-algebraic functions. Let  $S \subset R^k$  and  $T \subset R^\ell$  be semi-algebraic sets. A function  $f: S \rightarrow T$  is **semi-algebraic** if its graph  $\text{Graph}(f)$  is a semi-algebraic subset of  $R^{k+\ell}$ .

**Proposition 2.83.** *Let  $f: S \rightarrow T$  be a semi-algebraic function. If  $S' \subset S$  is semi-algebraic, then its image  $f(S')$  is semi-algebraic. If  $T' \subset T$  is semi-algebraic, then its inverse image  $f^{-1}(T')$  is semi-algebraic.*

**Proof:** The set  $f(S')$  is the image of  $(S' \times T) \cap \text{Graph}(f)$  under the projection from  $S \times T$  to  $T$  and is semi-algebraic by Theorem 2.76.

The set  $f^{-1}(T')$  is the image of  $(S \times T') \cap \text{Graph}(f)$  under the projection,  $S \times T \rightarrow S$  and is semi-algebraic, again by Theorem 2.76 □

**Proposition 2.84.** *If  $A, B, C$  are semi-algebraic sets in  $R^k, R^\ell$ , and  $R^m$ , resp., and  $f: A \rightarrow B$ ,  $g: B \rightarrow C$  are semi-algebraic functions, then the composite function  $g \circ f: A \rightarrow C$  is semi-algebraic.*

**Proof:** Let  $F \subset R^{k+\ell}$  be the graph of  $f$  and  $G \subset R^{\ell+m}$  the graph of  $g$ . The graph of  $g \circ f$  is the projection of  $(F \times R^m) \cap (R^k \times G)$  to  $R^{k+m}$  and hence is semi-algebraic by Theorem 2.76. □

**Proposition 2.85.** *Let  $A$  be a semi-algebraic set of  $\mathbb{R}^k$ . The semi-algebraic functions from  $A$  to  $\mathbb{R}$  form a ring.*

**Proof:** Follows from Proposition 2.84 by noting that  $f + g$  is the composition of  $(f, g): A \rightarrow \mathbb{R}^2$  with  $+: \mathbb{R}^2 \rightarrow \mathbb{R}$ , and  $f \times g$  is the composition of  $(f, g): A \rightarrow \mathbb{R}^2$  with  $\times: \mathbb{R}^2 \rightarrow \mathbb{R}$ .  $\square$

**Proposition 2.86.** *Let  $S \subset \mathbb{R}$  be a semi-algebraic set, and  $\varphi: S \rightarrow \mathbb{R}$  a semi-algebraic function. There exists a non-zero polynomial  $P \in \mathbb{R}[X, Y]$  such that for every  $x$  in  $S$ ,  $P(x, \varphi(x)) = 0$ .*

**Proof:** The graph  $\Gamma$  of  $\varphi$  is the finite union of non-empty semi-algebraic sets of the form

$$\Gamma_i = \{(x, y) \in \mathbb{R} \times \mathbb{R} \mid P_i(x, y) = 0 \wedge Q_{i,1}(x, y) > 0 \wedge \dots \wedge Q_{i,m_i}(x, y) > 0\}$$

with  $P_i$  not identically zero, for otherwise, given  $(x, y) \in \Gamma_i$ , the graph of  $\varphi$  intersected with the line  $X = x$  would contain a non-empty interval of this line. We can then take  $P$  as the product of the  $P_i$ .  $\square$

### 2.5.3 Extension of Semi-Algebraic Sets and Functions

In the following paragraphs,  $\mathbb{R}$  denotes a real closed field and  $\mathbb{R}'$  a real closed field containing  $\mathbb{R}$ . Given a semi-algebraic set  $S$  in  $\mathbb{R}^k$ , the **extension** of  $S$  to  $\mathbb{R}'$ , denoted  $\text{Ext}(S, \mathbb{R}')$ , is the semi-algebraic subset of  $\mathbb{R}'^k$  defined by the same quantifier free formula that defines  $S$ .

The following proposition is an easy consequence of Theorem 2.80.

**Proposition 2.87.** *Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set. The set  $\text{Ext}(S, \mathbb{R}')$  is well defined (i.e. it only depends on the set  $S$  and not on the quantifier free formula chosen to describe it).*

*The mapping  $S \rightarrow \text{Ext}(S, \mathbb{R}')$  preserves the boolean operations (finite intersection, finite union, and complementation).*

*If  $S \subset T$ , with  $T \subset \mathbb{R}^k$  semi-algebraic, then  $\text{Ext}(S, \mathbb{R}') \subset \text{Ext}(T, \mathbb{R}')$ .*

Of course  $\text{Ext}(S, \mathbb{R}') \cap \mathbb{R}^k = S$ . But  $\text{Ext}(S, \mathbb{R}')$  may not be the only semi-algebraic set of  $\mathbb{R}'^k$  with this property: if  $S = [0, 4] \subset \mathbb{R}_{\text{alg}}$  (the real algebraic numbers),  $\text{Ext}(S, \mathbb{R}) = [0, 4] \subset \mathbb{R}$ ; but also  $([0, \pi] \cup (\pi, 4]) \cap \mathbb{R}_{\text{alg}} = S$ , where  $\pi = 3.14\dots$  is the area enclosed by the unit circle.

**Exercise 2.16.** Show that if  $S$  is a finite semi-algebraic subset of  $\mathbb{R}^k$ , then  $\text{Ext}(S, \mathbb{R}')$  is equal to  $S$ .

For any real closed field  $\mathbb{R}$ , we denote by  $\pi$  the projection mapping

$$\pi: \mathbb{R}^{k+1} \rightarrow \mathbb{R}^k$$

that “forgets” the last coordinate.

**Proposition 2.88.** *If  $R$  is a real closed field and  $S \subset R^{k+1}$  is a semi-algebraic set then  $\pi(S)$  is semi-algebraic. Moreover, if  $R'$  is an arbitrary real closed extension of  $R$ , then  $\pi(\text{Ext}(S, R')) = \text{Ext}(\pi(S), R')$ .*

**Proof:** We use Theorem 2.80. Since the projection of the semi-algebraic set  $S$  is the semi-algebraic set  $B$ ,  $B = \pi(S)$  is true in  $R$ . This is expressed by a formula which is thus also true in  $R'$ . □

Let  $S \subset R^k$  and  $T \subset R^\ell$  be semi-algebraic sets, and let  $f: S \rightarrow T$  be a semi-algebraic function whose graph is  $G \subset S \times T$ .

**Proposition 2.89.** *If  $R'$  is a real closed extension of  $R$ , then  $\text{Ext}(G, R')$  is the graph of a semi-algebraic function  $\text{Ext}(f, R'): \text{Ext}(S, R') \rightarrow \text{Ext}(T, R')$ .*

**Proof:** Let  $\Phi, \Psi$  and  $\Gamma$  be quantifier free formulas such that

$$\begin{aligned} S &= \{x \in R^k \mid \Phi(x)\} \\ T &= \{y \in R^\ell \mid \Psi(y)\} \\ G &= \{(x, y) \in R^{k+\ell} \mid \Gamma(x, y)\}. \end{aligned}$$

The fact that  $G$  is the graph of a function from  $S$  to  $T$  can be expressed by the sentence  $\forall X A$ , with

$$\begin{aligned} A &= ((\Phi(X) \Leftrightarrow (\exists Y \Gamma(X, Y)) \wedge (\forall Y \Gamma(X, Y) \Rightarrow \Psi(Y)) \\ &\quad \wedge (\forall Y \forall Y' (\Gamma(X, Y) \wedge \Gamma(X, Y') \Rightarrow Y = Y'))), \end{aligned}$$

with  $X = (X_1, \dots, X_k)$ ,  $Y = (Y_1, \dots, Y_\ell)$  and  $Y' = (Y'_1, \dots, Y'_\ell)$ .

Applying Theorem 2.80,  $\forall X A$  is therefore true in  $R'$ , which expresses the fact that  $\text{Ext}(G, R')$  is the graph of a function from  $\text{Ext}(S, R')$  to  $\text{Ext}(T, R')$ , since

$$\begin{aligned} \text{Ext}(S, R') &= \{x \in R'^k \mid \Phi(x)\} \\ \text{Ext}(T, R') &= \{y \in R'^\ell \mid \Psi(y)\} \\ \text{Ext}(G, R') &= \{(x, y) \in R'^{k+\ell} \mid \Gamma(x, y)\}. \end{aligned}$$

□

The semi-algebraic function  $\text{Ext}(f, R')$  of the previous proposition is called the **extension of  $f$  to  $R'$** .

**Proposition 2.90.** *Let  $S'$  be a semi-algebraic subset of  $S$ . Then*

$$\text{Ext}(f(S'), R') = \text{Ext}(f, R')(\text{Ext}(S', R')).$$

**Proof:** The semi-algebraic set  $f(S')$  is the projection of  $G \cap (S' \times R^\ell)$  onto  $R^\ell$ , so the conclusion follows from Proposition 2.88. □



**Exercise 2.17.**

- a) Show that the semi-algebraic function  $f$  is injective (resp. surjective, resp. bijective) if and only if  $\text{Ext}(f, \mathbb{R}')$  is injective (resp. surjective, resp. bijective).
- b) Let  $T'$  be a semi-algebraic subset of  $T$ . Show that

$$\text{Ext}(f^{-1}(T'), \mathbb{R}') = \text{Ext}(f, \mathbb{R}')^{-1}(\text{Ext}(T', \mathbb{R}')) .$$

**2.6 Puiseux Series**

The field of Puiseux series provide an important example of a non-archimedean real closed field.

The collection of Puiseux series in  $\varepsilon$  with coefficients in  $\mathbb{R}$  will be a real closed field containing the field  $\mathbb{R}(\varepsilon)$  of rational functions in the variable  $\varepsilon$  ordered by  $0_+$  (see Notation 2.5). In order to include in our field roots of equations such as  $X^2 - \varepsilon = 0$ , we introduce rational exponents such as  $\varepsilon^{1/2}$ . This partially motivates the following definition of Puiseux series.

Let  $\mathbb{K}$  be a field and  $\varepsilon$  a variable. The **ring of formal power series in  $\varepsilon$  with coefficients in  $\mathbb{K}$** , denoted  $\mathbb{K}[[\varepsilon]]$ , consists of series of the form  $a = \sum_{i \geq 0} a_i \varepsilon^i$  with  $i \in \mathbb{N}$ ,  $a_i \in \mathbb{K}$ .

Its field of quotients, denoted  $\mathbb{K}((\varepsilon))$ , is the **field of Laurent series in  $\varepsilon$  with coefficients in  $\mathbb{K}$**  and consists of series of the form  $\bar{a} = \sum_{i \geq k} a_i \varepsilon^i$  with  $k \in \mathbb{Z}$ ,  $i \in \mathbb{Z}$ ,  $a_i \in \mathbb{K}$ .

**Exercise 2.18.** Prove that  $\mathbb{K}((\varepsilon))$  is a field, and is the quotient field of  $\mathbb{K}[[\varepsilon]]$ .

A **Puiseux series in  $\varepsilon$  with coefficients in  $\mathbb{K}$**  is a series of the form  $a = \sum_{i \geq k} a_i \varepsilon^{i/q}$  with  $k \in \mathbb{Z}$ ,  $i \in \mathbb{Z}$ ,  $a_i \in \mathbb{K}$ ,  $q$  a positive integer. Puiseux series are formal Laurent series in the indeterminate  $\varepsilon^{1/q}$  for some positive integer  $q$ . The **field of Puiseux series in  $\varepsilon$  with coefficients in  $\mathbb{K}$**  is denoted  $\mathbb{K}\langle\langle\varepsilon\rangle\rangle$ .

These series are formal in the sense that there is no assertion of convergence;  $\varepsilon$  is simply an indeterminate. We assume that the different symbols  $\varepsilon^r$ ,  $r \in \mathbb{Q}$ , satisfy

$$\begin{aligned} \varepsilon^{r_1} \varepsilon^{r_2} &= \varepsilon^{r_1+r_2}, \\ (\varepsilon^{r_1})^{r_2} &= \varepsilon^{r_1 r_2}, \\ \varepsilon^0 &= 1. \end{aligned}$$

Hence any two Puiseux series,  $\bar{a} = \sum_{i \geq k_1} a_i \varepsilon^{i/q_1}$ ,  $\bar{b} = \sum_{j \geq k_2} b_j \varepsilon^{j/q_2}$  can be written as formal Laurent series in  $\varepsilon^{1/q}$ , where  $q$  is the least common multiple of  $q_1$  and  $q_2$ . Thus, it is clear how to add and multiply two Puiseux series. Also, any finite number of Puiseux series can be written as formal Laurent series in  $\varepsilon^{1/q}$  with a common  $q$ .

If  $\bar{a} = a_1 \varepsilon^{r_1} + a_2 \varepsilon^{r_2} + \dots \in K\langle\langle\varepsilon\rangle\rangle$ , (with  $a_1 \neq 0$  and  $r_1 < r_2 < \dots$ ), then the **order** of  $\bar{a}$ , denoted  $o(\bar{a})$ , is  $r_1$  and the **initial coefficient** of  $\bar{a}$ , denoted  $\text{In}(\bar{a})$  is  $a_1$ . By convention, the order of 0 is  $\infty$ . The order is a function from  $K\langle\langle\varepsilon\rangle\rangle$  to  $\mathbb{Q} \cup \{\infty\}$  satisfying

- $o(\bar{a}\bar{b}) = o(\bar{a}) + o(\bar{b})$ ,
- $o(\bar{a} + \bar{b}) \geq \min(o(\bar{a}), o(\bar{b}))$ , with equality if  $o(\bar{a}) \neq o(\bar{b})$ .

**Exercise 2.19.** Prove that  $K\langle\langle\varepsilon\rangle\rangle$  is a field.

When  $K$  is an ordered field, we make  $K\langle\langle\varepsilon\rangle\rangle$  an ordered field by defining a Puiseux series  $\bar{a}$  to be positive if  $\text{In}(\bar{a})$  is positive. It is clear that the field of rational functions  $K(\varepsilon)$  equipped with the order  $0_+$  is a subfield of the ordered field of Puiseux series  $K\langle\langle\varepsilon\rangle\rangle$ , using Laurent’s expansions about 0.

In the ordered field  $K\langle\langle\varepsilon\rangle\rangle$ ,  $\varepsilon$  is infinitesimal over  $K$  (Definition page 32), since it is positive and smaller than any positive  $r \in K$ , since  $r - \varepsilon > 0$ . Hence, the field  $K\langle\langle\varepsilon\rangle\rangle$  is non-archimedean. This is the reason why we have chosen to name the indeterminate  $\varepsilon$  rather than some more neutral  $X$ .

The remainder of this section is primarily devoted to a proof of the following theorem.

**Theorem 2.91.** *Let  $R$  be a real closed field. Then, the field  $R\langle\langle\varepsilon\rangle\rangle$  is real closed.*

As a corollary

**Theorem 2.92.** *Let  $C$  be an algebraically closed field of characteristic 0. The field  $C\langle\langle\varepsilon\rangle\rangle$  is algebraically closed.*

**Proof:** Apply Theorem 2.31, Theorem 2.11 and Theorem 2.91, noticing that  $R[i]\langle\langle\varepsilon\rangle\rangle = R\langle\langle\varepsilon\rangle\rangle[i]$ . □

The first step in the proof of Theorem 2.91 is to show is that positive elements of  $R\langle\langle\varepsilon\rangle\rangle$  are squares in  $R\langle\langle\varepsilon\rangle\rangle$ .

**Lemma 2.93.** *A positive element of  $R\langle\langle\varepsilon\rangle\rangle$  is the square of an element in  $R\langle\langle\varepsilon\rangle\rangle$ .*

**Proof:** Suppose that  $\bar{a} = \sum_{i \geq k} a_i \varepsilon^{i/q} \in R\langle\langle\varepsilon\rangle\rangle$  with  $a_k > 0$ . Defining  $\bar{b} = \sum_{i \geq k+1} (a_i/a_k) \varepsilon^{(i-k)/q}$ , we have  $\bar{a} = a_k \varepsilon^{k/q} (1 + \bar{b})$  and  $o(\bar{b}) > 0$ .

The square root of  $1 + \bar{b}$  is obtained by taking the Taylor series expansion of  $(1 + \bar{b})^{1/2}$  which is

$$\bar{c} = 1 + \frac{1}{2} \bar{b} + \dots + \frac{1}{n!} \frac{1}{2} \left( \frac{1}{2} - 1 \right) \dots \left( \frac{1}{2} - (n-1) \right) \bar{b}^n + \dots$$

In order to check that  $\bar{c}^2 = 1 + \bar{b}$ , just substitute. Since  $a_k > 0$  and  $R$  is real closed,  $\sqrt{a_k} \in R$ . Hence,  $\sqrt{a_k} \varepsilon^{k/2q} \bar{c}$  is the square root of  $\bar{a}$ . □

In order to complete the proof of Theorem 2.91, it remains to prove that an odd degree polynomial in  $\mathbb{R}\langle\langle\varepsilon\rangle\rangle[X]$  has a root in  $\mathbb{R}\langle\langle\varepsilon\rangle\rangle$ . Given

$$P(X) = \bar{a}_0 + \bar{a}_1 X + \cdots + \bar{a}_p X^p \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle[X]$$

with  $p$  odd, we will construct an  $\bar{x} \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle$  such that  $P(\bar{x}) = 0$ . We may assume that  $\bar{a}_0 \neq 0$ , since otherwise 0 is a root of  $P$ . Furthermore, we may assume without loss of generality that

$$o(\bar{a}_i) = \frac{m_i}{m}$$

with the same  $m$  for every  $0 \leq i \leq p$ . Our strategy is to consider an unknown

$$\bar{x} = x_1 \varepsilon^{\xi_1} + x_2 \varepsilon^{\xi_1 + \xi_2} + \cdots + x_i \varepsilon^{\xi_1 + \cdots + \xi_i} + \cdots \quad (2.7)$$

with  $\xi_2 > 0, \dots, \xi_j > 0$  and determine, one after the other, the unknown coefficients  $x_i$  and the unknown exponents  $\xi_i$  so that  $\bar{x} \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle$  and satisfies  $P(\bar{x}) = 0$ .

Natural candidates for the choice of  $\xi_1$  and  $x_1$  will follow from the geometry of the exponents of  $P$ , that we study now. The polynomial  $P(X)$  can be thought of as a formal sum of expressions  $X^i \varepsilon^r$  ( $i \in \mathbb{Z}, r \in \mathbb{Q}$ ) with coefficients in  $\mathbb{R}$ . The points  $(i, r)$  for which  $X^i \varepsilon^r$  occurs in  $P(X)$  with non-zero coefficient constitute the **Newton diagram** of  $P$ . Notice that the points of the Newton diagram are arranged in columns and that the points  $M_i = (i, o(a_i))$ ,  $i = 0, \dots, p$ , for which  $\bar{a}_i \neq 0$  are the lowest points in each column.

The **Newton polygon** of  $P$  is the sequence of points

$$M_0 = M_{i_0}, \dots, M_{i_\ell} = M_p$$

satisfying:

- All points of the Newton diagram of  $P$  lie on or above each of the lines joining  $M_{i_{j-1}}$  to  $M_{i_j}$  for  $j = 1, \dots, \ell$ .
- The ordered triple of points  $M_{i_{j-1}}, M_{i_j}, M_{i_{j+1}}$  is oriented counter-clockwise, for  $j = 1, \dots, \ell - 1$ . This is saying that the edges joining adjacent points in the sequence  $M_0 = M_{i_0}, \dots, M_{i_\ell} = M_p$  constitute a **convex chain**.

In such a case the **slope** of  $[M_{i_{j-1}}, M_{i_j}]$  is  $\frac{o(a_{i_j}) - o(a_{i_{j-1}})}{i_j - i_{j-1}}$ , and its **horizontal projection** is the interval  $[i_{j-1}, i_j]$ .

Notice that the Newton polygon of  $P$  is the lower convex hull of the Newton diagram of  $P$ .

To the segment  $E = [M_{i_{j-1}}, M_{i_j}]$  with horizontal projection  $[i_{j-1}, i_j]$ , we associate its **characteristic polynomial**

$$Q(P, E, X) = \sum a_h X^h \in \mathbb{R}[X],$$

where the sum is over all  $h$  for which

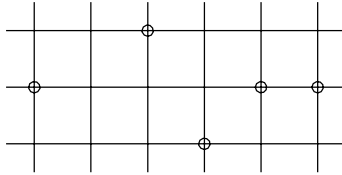
$$M_h = (h, o(\bar{a}_h)) = \left( h, \frac{m_h}{m} \right) \in E \text{ and } a_h = \text{In}(\bar{a}_h).$$

Note that if  $-\xi$  is the slope of  $E$ , then  $o(\bar{a}_h) + h\xi$  has a constant value  $\beta$  for all  $M_h$  on  $E$ .

*Example 2.94.* Let

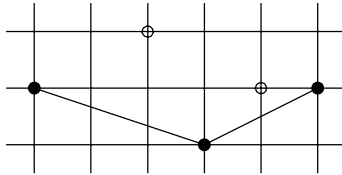
$$P(X) = \varepsilon - 2\varepsilon^2 X^2 - X^3 + \varepsilon X^4 + \varepsilon X^5.$$

The Newton diagram of  $P$  is



**Fig. 2.3.** Newton diagram

The Newton polygon of  $P$  consists of two segments  $E = [M_0, M_3]$  and  $F = [M_3, M_5]$ . The segment  $E$  has an horizontal projection of length 3 and the segment  $F$  has an horizontal projection of length 2



**Fig. 2.4.** Newton polygon

We have

$$\begin{aligned} Q(P, E, X) &= 1 - X^3 \\ Q(P, F, X) &= X^3(X^2 - 1). \end{aligned}$$

The two slopes are  $-1/3$  and  $1/2$  and the corresponding values of  $\xi$  are  $1/3$  and  $-1/2$ . The common value  $\beta$  of  $o(\bar{a}_h) + h\xi$  on the two segments are 1 and  $-3/2$ . □

If  $x$  is a non-zero root of multiplicity  $r$  of the characteristic polynomial of a segment  $E$  of the Newton polygon with slope  $-\xi$ , we construct a root of  $P$  which is a Puiseux series starting with  $x\varepsilon^\xi$ . In other words we find

$$\bar{x} = x\varepsilon^\xi + x_2\varepsilon^{\xi+\xi_2} + \dots + x_i\varepsilon^{\xi+\xi_2+\dots+\xi_i} + \dots \tag{2.8}$$

with  $\xi_2 > 0, \dots, \xi_j > 0$  such that  $P(\bar{x}) = 0$ .

The next lemma is a key step in this direction. The result is the following: if we replace in  $P X$  by  $\varepsilon^\xi(x + X)$  and divide the result by  $\varepsilon^{-\beta}$ , where  $\beta$  is the common value of  $o(\bar{a}_h) + h\xi$  on  $E$ , we obtain a new Newton polygon with a part having only negative slopes, whose horizontal projection is  $[0, r]$ . A segment of this part of the Newton polygon will be used to find the second term of the series.

**Lemma 2.95.** *Let*

- $\xi$  be the opposite of the slope of a segment  $E$  of the Newton polygon of  $P$ ,
- $\beta$  be the common value of  $o(\bar{a}_h) + h\xi$  for all  $q_h$  on  $E$ ,
- $x \in \mathbb{R}$  be a non-zero root of the characteristic polynomial  $Q(P, E, X)$  of multiplicity  $r$ .

a) *The polynomial*

$$R(P, E, x, Y) = \varepsilon^{-\beta} P(\varepsilon^\xi(x + Y)) = \bar{b}_0 + \bar{b}_1 Y + \dots + \bar{b}_p Y^p$$

satisfies

$$\begin{aligned} o(\bar{b}_i) &\geq 0, & i = 0, \dots, p, \\ o(\bar{b}_i) &> 0, & i = 0, \dots, r - 1, \\ o(\bar{b}_r) &= 0. \end{aligned}$$

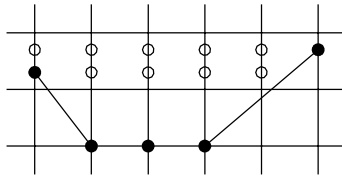
b) *For every  $\bar{x} \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle$  such that  $\bar{x} = \varepsilon^\xi(x + \bar{y})$  with  $o(\bar{y}) > 0$ ,  $o(P(\bar{x})) > \beta$ .*

We illustrate the construction in our example.

*Example 2.96.* Continuing Example 2.94, we choose the segment  $E$ , with  $\xi = 1/3$ , chose the root  $x = 1$  of  $X^3 - 1$ , with multiplicity 1, and replace  $X$  by  $\varepsilon^{1/3}(1 + X)$  and get

$$\begin{aligned} P_1(X) &= \varepsilon^{-1} P(\varepsilon^{1/3}(1 + X)) \\ &= \varepsilon^{5/3} X^5 + \left(\varepsilon^{4/3} + 5\varepsilon^{5/3}\right) X^4 \\ &\quad + \left(-1 + 4\varepsilon^{4/3} + 10\varepsilon^{5/3}\right) X^3 \\ &\quad + \left(-3 + 8\varepsilon^{5/3} + 6\varepsilon^{4/3}\right) X^2 \\ &\quad + \left(\varepsilon^{5/3} - 3 + 4\varepsilon^{4/3}\right) X - \varepsilon^{5/3} + \varepsilon^{4/3}. \end{aligned}$$

The Newton polygon of  $p_1$  is



**Fig. 2.5.** Newton polygon of  $p_1$

We chose the negative slope with corresponding characteristic polynomial  $-3X + 1$  and make the change of variable  $X = \varepsilon^{4/3}(1/3 + Y)$ .

We have obtained this way the two first terms  $\varepsilon^{1/3} + (1/3)\varepsilon^{1/3+4/3} + \dots$  of a Puiseux series  $\bar{x}$  satisfying  $P(\bar{x}) = 0$ . □

The proof of Lemma 2.95 uses the next lemma which describes a property of the characteristic polynomials associated to the segments of the Newton polygon.

**Lemma 2.97.** *The slope  $-\xi$  of  $E$  has the form  $-c/(mq)$  with  $q > 0$  and  $\gcd(c, q) = 1$ . Moreover,  $Q(P, E, X) = X^j \phi(X^q)$ , where  $\phi \in \mathbb{R}[X]$ ,  $\phi(0) \neq 0$ , and  $\deg \phi = (k - j)/q$ .*

**Proof:** The slope of  $E = [M_j, M_k]$  is

$$\frac{o(\bar{a}_k) - o(\bar{a}_j)}{k - j} = \frac{m_k - m_j}{m(k - j)} = -\frac{c}{mq},$$

where  $q > 0$  and  $\gcd(c, q) = 1$ . If  $(h, o(\bar{a}_h)) = (h, \frac{mh}{m})$  is on  $E$  then

$$\frac{c}{mq} = \frac{o(\bar{a}_j) - o(\bar{a}_h)}{h - j} = \frac{m_j - m_h}{m(h - j)}.$$

Hence,  $q$  divides  $h - j$ , and there exists a non-negative  $s$  such that  $h = j + sq$ . The claimed form of  $Q(P, E, X)$  follows. □

**Proof of Lemma 2.95:** For a) since  $x$  is a root of  $\phi(X^q)$  of multiplicity  $r$ , we have

$$\phi(X^q) = (X - x)^r \psi(X), \quad \psi(x) \neq 0.$$

Thus,

$$\begin{aligned} R(P, E, x, Y) &= \varepsilon^{-\beta} P(\varepsilon^\xi(x + Y)) \\ &= \varepsilon^{-\beta} (\bar{a}_0 + \bar{a}_1 \varepsilon^\xi(x + Y) + \dots + \bar{a}_p \varepsilon^{p\xi}(x + Y)^p) \\ &= A(Y) + B(Y), \end{aligned}$$

where

$$\begin{aligned} A(Y) &= \varepsilon^{-\beta} \sum_{h, qh \in E} a_h \varepsilon^{o(a_h) + h\xi} (x + Y)^h \\ B(Y) &= \varepsilon^{-\beta} \left( \sum_{h, qh \in E} (\bar{a}_h - a_h \varepsilon^{o(a_h)}) \varepsilon^{h\xi} (x + Y)^h + \sum_{\ell, q\ell \notin E} \bar{a}_\ell \varepsilon^{\ell\xi} (x + Y)^\ell \right). \end{aligned}$$

Since  $o(\bar{a}_h) + h\xi = \beta$ ,

$$\begin{aligned} A(Y) &= Q(P, E, x + Y) \\ &= (x + Y)^j \phi((x + Y)^q) \\ &= Y^r (x + Y)^j \psi(x + Y) \\ &= c_r Y^r + c_{r+1} Y^{r+1} + \dots + c_p Y^p, \end{aligned}$$

with  $c_r = x^j \psi(x) \neq 0$  and  $c_i \in \mathbb{R}$ .

Since  $o((\bar{a}_h - a_h \varepsilon^{o(a_h)}) \varepsilon^{h\xi}) > \beta$  and  $o(\bar{a}_\ell \varepsilon^{\ell\xi}) > \beta$ ,

$$R(P, E, x, Y) = B(Y) + c_r Y^r + c_{r+1} Y^{r+1} + \dots + c_p Y^p,$$

where every coefficient of  $B(Y) \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle[Y]$  has positive order. The conclusion follows.

For b), since  $o(\bar{y}) > 0$ ,  $o(R(P, E, x, \bar{y})) > 0$  is an easy consequence of a). The conclusion follows noting that  $P(x) = \varepsilon^\beta R(P, E, x, y)$ .  $\square$

It is now possible to proceed with the proof of Theorem 2.91.

**Proof of Theorem 2.91:** Consider  $P$  with odd degree. Hence, we can choose a segment  $E_1$  of the Newton polygon of  $P$  which has a horizontal projection of odd length. Let the slope of  $E_1$  be  $-\xi_1$ . It follows from Lemma 2.97 that the corresponding characteristic polynomial  $Q(P, E_1, X)$  has a non-zero root  $x_1$  in  $\mathbb{R}$  of odd multiplicity  $r_1$ , since  $\mathbb{R}$  is real closed. Define  $P_1(X) = R(P, E_1, x_1, X)$  using this segment and the root  $x_1$ .

Note that  $(r_1, 0)$  is a vertex of the Newton polygon of  $R(P, E_1, x_1, X)$ , and that all the slopes of segments  $[M_j, M_k]$  of the Newton polygon of  $R(P, E_1, x_1, X)$  for  $k \leq r_1$  are negative: this is an immediate consequence of Lemma 2.95.

Choose recursively a segment  $E_{i+1}$  of the Newton polygon of  $P_i$  with negative slope  $-\xi_{i+1}$ , and horizontal projection of odd length, so that the corresponding characteristic polynomial  $Q(P_i, E_{i+1}, X)$  has a non-zero root  $x_{i+1}$  in  $\mathbb{R}$  of odd multiplicity  $r_{i+1}$ , and take  $P_{i+1}(X) = R(P_i, E_{i+1}, x_{i+1}, X)$ . The only barrier to continuing this process is if we cannot choose a segment with negative slope over the interval  $[0, r_i]$  and this is the case only if 0 is a root of  $P_i(X)$ . But in this exceptional case  $x_1 \varepsilon^{\xi_1} + \dots + x_i \varepsilon^{\xi_1 + \dots + \xi_i}$  is clearly a root of  $P$ .

Suppose we have constructed  $x_i, \xi_i$  for  $i \in \mathbb{N}$  and let

$$\bar{x} = x_1 \varepsilon^{\xi_1} + x_2 \varepsilon^{\xi_1 + \xi_2} + \dots$$

Then from the definition of the  $P_i(X)$ , it follows by induction that  $o(P(x)) > \beta_1 + \dots + \beta_j$  for all  $j$ . To complete the proof, we need to know that  $\bar{x} \in \mathbb{R}\langle\langle\varepsilon\rangle\rangle$  and that the sums  $\beta_1 + \dots + \beta_j$  are unbounded. Both these will follow if we know that the  $q$  in Lemma 2.97 is eventually 1. Note that the multiplicities of the chosen roots  $x_i$  are non-increasing and hence are eventually constant, at which point they have the value  $r$ . This means that from this point on, the Newton polygon has a single segment with negative slope, and horizontal projection of length  $r$ . Therefore all subsequent roots chosen also have multiplicity  $r$ . It follows (since  $Q_j(X)$  must also have degree  $r$ ) that  $Q_j(X) = c(X - x_j)^r$  with  $x_j \neq 0$ , from which it follows that the corresponding  $q$  is equal to 1, since the coefficient of degree 1 of  $\phi_j$  is  $-rcx_j^{r-1}$ , which is not zero.  $\square$

If  $K$  is a field, we denote by  $K\langle\varepsilon\rangle$  the subfield of  $K\langle\langle\varepsilon\rangle\rangle$  of **algebraic Puiseux series**, which consists of those elements that are algebraic over  $K(\varepsilon)$ , i.e. that satisfy a polynomial equation with coefficients in  $K(\varepsilon)$ .

**Corollary 2.98.** *When  $R$  is real closed,  $R\langle\varepsilon\rangle$  is real closed. The field  $R\langle\varepsilon\rangle$  is the real closure of  $R(\varepsilon)$  equipped with the order  $0_+$ .*

**Proof:** Follows immediately from Theorem 2.91 and Exercise 2.10.  $\square$

Similarly, if  $C = R[i]$ , then  $C\langle\varepsilon\rangle = R\langle\varepsilon\rangle[i]$  is an algebraic closure of  $C(\varepsilon)$ .

We shall see in Chapter 3 that algebraic Puiseux series with coefficients in  $R$  can be interpreted as of germs semi-algebraic and continuous functions at the right of the origin.

A **valuation ring** of a field  $F$  is a subring of  $F$  such that either  $x$  or its inverse is in the ring for every non-zero  $x$ .

**Proposition 2.99.** *The elements of  $K\langle\varepsilon\rangle$  with non-negative order constitute a valuation ring denoted  $K\langle\varepsilon\rangle_b$ . The elements of  $R\langle\varepsilon\rangle_b$  are exactly the elements of  $R\langle\varepsilon\rangle$  bounded over  $R$  (i.e. their absolute value is less than a positive element of  $R$ ). The elements of  $C\langle\varepsilon\rangle_b$  are exactly the elements of  $C\langle\varepsilon\rangle$  bounded over  $R$  (i.e. their modulus is less than a positive element of  $R$ ).*

**Notation 2.100. [Limit]** We denote by  $\lim_\varepsilon$  the ring homomorphism from  $K\langle\varepsilon\rangle_b$  to  $K$  which maps  $\sum_{i \in \mathbb{N}} a_i \varepsilon^{i/q}$  to  $a_0$ . The mapping  $\lim_\varepsilon$  simply replaces  $\varepsilon$  by 0 in a bounded Puiseux series.  $\square$

## 2.7 Bibliographical Notes

The theory of real closed fields was developed by Artin and Schreier [7] and used by Artin [6] in his solution to Hilbert's 17-th problem. The algebraic proof of the fundamental theorem of algebra is due to Gauss [65].

Real root counting began with Descartes's law of sign [53], generalized by Budan [34] and Fourier [60], and continued with Sturm [152]. The connection between virtual roots [68] and Budan-Fourier's theorem comes from [49]. The notion of Cauchy index appears in [41]. Theorem 2.58 is already proved in two particular cases (when  $Q = P'$  and when  $P$  is square-free) in [152]. The partial converse to Descartes's law of sign presented here appears in [126].

Quantifier elimination for real closed fields is a fundamental result. It was known to Tarski before 1940 (it is announced without a proof in [154]) and published much later [156]. The version of 1940, ready for publication in *Actualités Scientifiques et Industrielles* (Hermann), was finally not published at that time, "as a result of war activities", and has appeared in print much later [155]. The proof presented here follows the original procedure of Tarski. Theorem 2.61 is explicitly stated in [155, 156], and the sign determination algorithm is sketched.



There are many different proofs of quantifier elimination for real closed fields, in particular by Seidenberg [148], Cohen [43] and Hormander [92].

Puiseux series have been considered for the first time by Newton [123].

---

## Semi-Algebraic Sets

In Section 3.1 and Section 3.2, we define the topology of semi-algebraic sets and study connectedness in a general real closed field. In order to study the properties of closed and bounded semi-algebraic sets in Section 3.4, we introduce semi-algebraic germs in Section 3.3. The semi-algebraic germs over a real closed field constitute a real closed field containing infinitesimal elements, closely related to the field of Puiseux series seen in Chapter 2, and play an important role throughout the whole book. We end the chapter with Section 3.5 on semi-algebraic differentiable functions.

### 3.1 Topology

Let  $\mathbf{R}$  be a real closed field. Since  $\mathbf{R}$  is an ordered field, we can define the topology on  $\mathbf{R}^k$  in terms of open balls in essentially the same way that we define the topology on  $\mathbb{R}^k$ . The **euclidean norm**, **open balls**, **closed balls**, and **spheres** are defined as follows:

With  $x = (x_1, \dots, x_k) \in \mathbf{R}^k$ ,  $r \in \mathbf{R}$ ,  $r > 0$ , we denote

$$\begin{aligned} \|x\| &= \sqrt{x_1^2 + \dots + x_k^2} && \text{(euclidean norm of } x), \\ B_k(x, r) &= \{y \in \mathbf{R}^k \mid \|y - x\|^2 < r^2\} && \text{(open ball),} \\ \bar{B}_k(x, r) &= \{y \in \mathbf{R}^k \mid \|y - x\|^2 \leq r^2\} && \text{(closed ball),} \\ S^{k-1}(x, r) &= \{y \in \mathbf{R}^k \mid \|y - x\|^2 = r^2\} && \text{((} k - 1 \text{)-sphere).} \end{aligned}$$

Note that  $B_k(x, r)$ ,  $\bar{B}_k(x, r)$ , and  $S^{k-1}(x, r)$  are semi-algebraic sets.

We omit both  $x$  and  $r$  from the notation when  $x$  is the origin of  $\mathbf{R}^k$  and  $r = 1$ , i.e. for the unit ball and sphere centered at the origin. We also omit the subscript  $k$  when it leads to no ambiguity.

We recall the definitions of the basic notions of open, closed, closure, interior, continuity, etc.

A set  $U \subset \mathbb{R}^k$  is **open** if it is the union of open balls, i.e. if every point of  $U$  is contained in an open ball contained in  $U$ . A set  $F \subset \mathbb{R}^k$  is **closed** if its complement is open. Clearly, the arbitrary union of open sets is open and the arbitrary intersection of closed sets is closed. The **closure** of a set  $S$ , denoted  $\bar{S}$ , is the intersection of all closed sets containing  $S$ . The **interior** of  $S$ , denoted  $S^\circ$ , is the union of all open subsets of  $S$  and thus is also the union of all open balls in  $S$ . We also have a notion of subsets of  $S$  being open or closed relative to  $S$ . A subset of  $S$  is called **open in  $S$**  if it is the intersection of an open set with  $S$ . It is **closed in  $S$**  if it is the intersection of a closed set with  $S$ . A function from  $S$  to  $T$  is **continuous** if the inverse image of any set open in  $T$  is open in  $S$ . It is easy to prove that polynomial maps from  $\mathbb{R}^k$  to  $\mathbb{R}^\ell$  are continuous in the Euclidean topology: one proves first that  $+$  and  $\times$  are continuous, then that the composite of continuous functions is continuous.

These definitions are clearly equivalent to the following formulations:

- $U$  is open if and only if  $\forall x \in U \exists r \in \mathbb{R}, r > 0 B(x, r) \subset U$ .
- $\bar{S} = \{x \in \mathbb{R}^k \mid \forall r > 0 \exists y \in S \ \|y - x\|^2 < r^2\}$ .
- $S^\circ = \{x \in S \mid \exists r > 0, \forall y \ \|y - x\|^2 < r^2 \Rightarrow y \in S\}$ .
- If  $S \subset \mathbb{R}^k$  and  $T \subset \mathbb{R}^\ell$ , a function  $f: S \rightarrow T$  is continuous if and only if it is continuous at every point of  $S$ , i.e.

$$\forall x \in S \forall r > 0 \exists \delta > 0, \forall y \in S \ \|y - x\| < \delta \Rightarrow \|f(y) - f(x)\| < r.$$

Note that if  $U, S, T, f$  are semi-algebraic, these definitions are expressed by formulas in the language of ordered fields. Indeed, it is possible to replace in these definitions semi-algebraic sets and semi-algebraic functions by quantifier-free formulas describing them. For example let  $\Psi(X_1, \dots, X_k)$  be a quantifier free formula such that

$$S = \{(x_1, \dots, x_k) \in \mathbb{R}^k \mid \Psi(x_1, \dots, x_k)\}.$$

Then, if  $\Phi(X_1, \dots, X_k, Y_1, \dots, Y_\ell)$  is a formula,  $\forall x \in S \ \Phi(x, y)$  can be replaced by

$$(\forall x_1) \dots (\forall x_k) (\Psi(x_1, \dots, x_k) \Rightarrow \Phi(x_1, \dots, x_k, y_1, \dots, y_\ell)),$$

and  $\exists x \in S \ \Phi(x, y_1, \dots, y_\ell)$  can be replaced by

$$(\exists x_1) \dots (\exists x_k) (\Psi(x_1, \dots, x_k) \wedge \Phi(x_1, \dots, x_k, y_1, \dots, y_\ell)).$$

An immediate consequence of these observations and of Theorem 2.77 (Quantifier elimination) (more precisely Corollary 2.78) is

**Proposition 3.1.** *The closure and the interior of a semi-algebraic set are semi-algebraic sets.*

*Remark 3.2.* It is tempting to think that the closure of a semi-algebraic set is obtained by relaxing the strict inequalities describing the set, but this idea is mistaken. Take  $S = \{x \in \mathbb{R} \mid x^3 - x^2 > 0\}$ . The closure of  $S$  is not  $T = \{x \in \mathbb{R} \mid x^3 - x^2 \geq 0\}$  but is  $\bar{S} = \{x \in \mathbb{R} \mid x^3 - x^2 \geq 0 \wedge x \geq 1\}$ , as 0 is clearly in  $T$  and not in  $\bar{S}$ .  $\square$

We next consider semi-algebraic and continuous functions. The following proposition is clear, noting that Proposition 2.85 and Proposition 2.84 take care of the semi-algebraicity:

**Proposition 3.3.** *If  $A, B, C$  are semi-algebraic sets and  $f: A \rightarrow B$  and  $g: B \rightarrow C$  are semi-algebraic continuous functions, then the composite function  $g \circ f: A \rightarrow C$  is semi-algebraic and continuous.*

*Let  $A$  be a semi-algebraic set of  $\mathbb{R}^k$ . The semi-algebraic continuous functions from  $A$  to  $\mathbb{R}$  form a ring.*

**Exercise 3.1.** Let  $\mathbb{R}'$  be a real closed field containing  $\mathbb{R}$ .

- a) Show that the semi-algebraic set  $S \subset \mathbb{R}^k$  is open (resp. closed) if and only if  $\text{Ext}(S, \mathbb{R}')$  is open (resp. closed). Show that

$$\text{Ext}(\overline{S}, \mathbb{R}') = \overline{\text{Ext}(S, \mathbb{R}')}$$

- b) Show that a semi-algebraic function  $f$  is continuous if and only if  $\text{Ext}(f, \mathbb{R}')$  is continuous.

The intermediate value property is valid for semi-algebraic continuous functions.

**Proposition 3.4.** *Let  $f$  be a semi-algebraic and continuous function defined on  $[a, b]$ . If  $f(a)f(b) < 0$ , then there exists  $x$  in  $(a, b)$  such that  $f(x) = 0$ .*

**Proof:** Suppose, without loss of generality, that  $f(a) > 0, f(b) < 0$ . Let  $A = \{x \in [a, b] \mid f(x) > 0\}$ . The set  $A$  is semi-algebraic, non-empty, and open. So, by Corollary 2.79,  $A$  is the union of a finite non-zero number of open subintervals of  $[a, b]$ . Let  $A = [a, b_1) \cup \dots \cup (a_\ell, b_\ell]$ . Then  $f(b_1) = 0$  since  $f$  is continuous, thus  $f(b_1) \leq 0$ . □

**Proposition 3.5.** *Let  $f$  be a semi-algebraic function defined on the semi-algebraic set  $S$ . Then  $f$  is continuous if and only if for every  $x \in S$  and every  $y \in \text{Ext}(S, \mathbb{R}(\varepsilon))$  such that  $\lim_\varepsilon(y) = x, \lim_\varepsilon(\text{Ext}(f, \mathbb{R}(\varepsilon))(y)) = f(x)$ .*

**Proof:** Suppose that  $f$  is continuous. Then

$$\forall x \in S \forall a > 0 \exists b(a) \forall y \in S \mid x - y \mid < b(a) \Rightarrow \mid f(x) - f(y) \mid < a.$$

holds in  $\mathbb{R}$ . Taking  $y \in \text{Ext}(S, \mathbb{R}(\varepsilon))$  such that  $\lim_\varepsilon(y) = x$ , for every positive  $a \in \mathbb{R}, \mid x - y \mid < b(a)$ , thus  $\mid f(x) - \text{Ext}(f, \mathbb{R}(\varepsilon))(y) \mid < a$ , using Tarski-Seidenberg principle (Theorem 2.80).

In the other direction, suppose that  $f$  is not continuous. It means that

$$\exists x \in S \exists a > 0 \forall b \exists y \in S \mid x - y \mid < b \wedge \mid f(x) - f(y) \mid > a$$

holds in  $\mathbb{R}$ . Taking  $b = \varepsilon$ , there exists  $y \in \text{Ext}(S, \mathbb{R}(\varepsilon))$  such that  $\lim_{\varepsilon} (y) = x$ , while  $|f(x) - \text{Ext}(f, \mathbb{R}(\varepsilon))(y)| > a$ , using again Tarski-Seidenberg principle (Theorem 2.80), which implies that  $f(x)$  and  $\lim_{\varepsilon} (\text{Ext}(f, \mathbb{R}(\varepsilon))(y))$  are not infinitesimally close. □

A **semi-algebraic homeomorphism**  $f$  from a semi-algebraic set  $S$  to a semi-algebraic set  $T$  is a semi-algebraic bijection which is continuous and such that  $f^{-1}$  is continuous.

**Exercise 3.2.** Let  $\mathbb{R}'$  be a real closed field containing  $\mathbb{R}$ . Prove that if  $f$  is a semi-algebraic homeomorphism from a semi-algebraic set  $S$  to a semi-algebraic set  $T$ , then  $\text{Ext}(f, \mathbb{R}')$  is a semi-algebraic homeomorphism from  $\text{Ext}(S, \mathbb{R}')$  to  $\text{Ext}(T, \mathbb{R}')$ .

### 3.2 Semi-algebraically Connected Sets

Recall that a set  $S \subset \mathbb{R}^k$  is connected if  $S$  is not the disjoint union of two non-empty sets which are both closed in  $S$ . Equivalently,  $S$  does not contain a non-empty strict subset which is both open and closed in  $S$ .

Unfortunately, this definition is too general to be suitable for  $\mathbb{R}^k$  with  $\mathbb{R}$  an arbitrary real closed field, as it allows  $\mathbb{R}$  to be disconnected.

For example, consider  $\mathbb{R}_{\text{alg}}$ , the field of real algebraic numbers. The set  $(-\infty, \pi) \cap \mathbb{R}_{\text{alg}}$  is both open and closed (with  $\pi = 3.14\dots$ ), and hence  $\mathbb{R}_{\text{alg}}$  is not connected. However, the set  $(-\infty, \pi) \cap \mathbb{R}_{\text{alg}}$  is not a semi-algebraic set in  $\mathbb{R}_{\text{alg}}$ , since  $\pi$  is not an algebraic number.

Since semi-algebraic sets are the only sets in which we are interested, we restrict our attention to these sets.

A semi-algebraic set  $S \subset \mathbb{R}^k$  is **semi-algebraically connected** if  $S$  is not the disjoint union of two non-empty semi-algebraic sets that are both closed in  $S$ . Or, equivalently,  $S$  does not contain a non-empty semi-algebraic strict subset which is both open and closed in  $S$ .

A semi-algebraic set  $S$  in  $\mathbb{R}^k$  is **semi-algebraically path connected** when for every  $x, y$  in  $S$ , there exists a **semi-algebraic path** from  $x$  to  $y$ , i.e. a continuous semi-algebraic function  $\varphi: [0, 1] \rightarrow S$  such that  $\varphi(0) = x$  and  $\varphi(1) = y$ .

We shall see later, in Chapter 5 (Theorem 5.23), that the two notions of being semi-algebraically connected and semi-algebraically path connected agree for semi-algebraic sets. We shall see also (Theorem 5.22) that the two notions of being connected and semi-algebraically connected agree for semi-algebraic subsets of  $\mathbb{R}$ .

**Exercise 3.3.** Prove that if  $A$  is semi-algebraically connected, and the semi-algebraic set  $B$  is semi-algebraically homeomorphic to  $A$  then  $B$  is semi-algebraically connected.

Since the semi-algebraic subsets of the real closed field  $\mathbb{R}$  are the finite unions of open intervals and points, the following proposition is clear:

**Proposition 3.6.** *A real closed field  $\mathbb{R}$  (as well as all its intervals) is semi-algebraically connected.*

A subset  $C$  of  $\mathbb{R}^k$  is **convex** if  $x, y \in C$  implies that the segment

$$[x, y] = \{(1 - \lambda)x + \lambda y \mid \lambda \in [0, 1] \subset \mathbb{R}\}$$

is contained in  $C$ .

**Proposition 3.7.** *If  $C$  is semi-algebraic and convex then  $C$  is semi-algebraically connected.*

**Proof:** Suppose that  $C$  is the disjoint union of two non-empty sets  $F_1$  and  $F_2$  which are closed in  $C$ . Let  $x_1 \in F_1$  and  $x_2 \in F_2$ . The segment  $[x_1, x_2]$  is the disjoint union of  $F_1 \cap [x_1, x_2]$  and  $F_2 \cap [x_1, x_2]$ , which are closed, semi-algebraic, and non-empty. This contradicts the fact that  $[x_1, x_2]$  is semi-algebraically connected (Proposition 3.6).  $\square$

Since the open cube  $(0, 1)^k$  is convex, the following proposition is clear:

**Proposition 3.8.** *The open cube  $(0, 1)^k$  is semi-algebraically connected.*

The following useful property holds for semi-algebraically connected sets.

**Proposition 3.9.** *If  $S$  is a semi-algebraically connected semi-algebraic set and  $f: S \rightarrow \mathbb{R}$  is a locally constant semi-algebraic function (i.e. given  $x \in S$ , there is an open  $U \subset S$  such that for all  $y \in U$ ,  $f(y) = f(x)$ ), then  $f$  is a constant.*

**Proof:** Let  $d \in f(S)$ . Since  $f$  is locally constant  $f^{-1}(d)$  is open. If  $f$  is not constant,  $f(S) \setminus \{d\}$  is non-empty and  $f^{-1}(f(S) \setminus \{d\})$  is open. Clearly,  $S = f^{-1}(d) \cup f^{-1}(f(S) \setminus \{d\})$ . This contradicts the fact that  $S$  is semi-algebraically connected, since  $f^{-1}(d)$  and  $f^{-1}(f(S) \setminus \{d\})$  are non-empty open and disjoint semi-algebraic sets.  $\square$

### 3.3 Semi-algebraic Germs

We introduce the field of germs of semi-algebraic continuous functions at the right of the origin and prove that it provides another description of the real closure  $\mathbb{R}(\varepsilon)$  of  $\mathbb{R}(\varepsilon)$  equipped with the order  $0_+$ . We saw in Chapter 2 that  $\mathbb{R}(\varepsilon)$  is the field of algebraic Puiseux series (Corollary 2.98). The field  $\mathbb{R}(\varepsilon)$  is used in Section 3.4 to prove results in semi-algebraic geometry, and it will also play an important role in the second part of the book, which is devoted to algorithms.

In order to define the field of germs of semi-algebraic continuous functions at the right of the origin, some preliminary work on semi-algebraic and continuous functions is necessary.

**Proposition 3.10.** *Let  $S$  be a semi-algebraic set and let  $P$  be a univariate polynomial with coefficients semi-algebraic continuous functions defined on  $S$ . Then if  $y$  is a simple root of  $P(x, Y)$  for a given  $x \in S$ , there is a semi-algebraic and continuous function  $f$  defined on a neighborhood of  $x$  in  $S$  such that  $f(x) = y$  and for every  $x' \in U$ ,  $f(x')$  is a simple root of  $P(x', Y)$ .*

**Proof:** Let  $m > 0$  such that for every  $m' \in (0, m)$ ,

$$P(x, y - m')P(x, y + m') < 0.$$

Such an  $m$  exists because,  $y$  being a simple root of  $P(x, Y)$ ,  $P(x, Y)$  is either increasing or decreasing on an interval  $(y - m, y + m)$ . Note that  $y$  is the only root of  $P(x, Y)$  in  $(y - m, y + m)$ . Suppose without loss of generality, that  $\partial P / \partial Y(x, y) > 0$  and let  $V$  be a neighborhood of  $(x, y)$  in  $S \times \mathbb{R}$  where  $\partial P / \partial Y$  is positive. For every  $m'$ ,  $0 < m' < m$ , the set

$$\{u \in S \mid P(u, y - m')P(u, y + m') < 0 \wedge [(u, y - m'), (u, y + m')] \subset V\}$$

is an open semi-algebraic subset of  $S$  containing  $x$ . This proves that  $P(u, Y)$  has a simple root  $y(u)$  on  $(y - m', y + m')$  and that the function associating to  $u \in U$  the value  $y(u)$  is continuous.  $\square$

The set of **germs of semi-algebraic continuous functions at the right of the origin** is the set of semi-algebraic continuous functions with values in  $\mathbb{R}$  which are defined on an interval of the form  $(0, t)$ ,  $t \in \mathbb{R}_+$ , modulo the equivalence relation

$$f_1 \simeq f_2 \Leftrightarrow \exists t > 0 \quad \forall t' \quad 0 < t' < t \quad f_1(t') = f_2(t').$$

**Proposition 3.11.** *The germs of semi-algebraic continuous functions at the right of the origin form a real closed field.*

**Proof:** Let  $\varphi$  and  $\varphi'$  be two germs of semi-algebraic continuous functions at the right of the origin, and consider semi-algebraic continuous functions  $f$  and  $f'$  representing  $\varphi$  and  $\varphi'$ , defined without loss of generality on a common interval  $(0, t)$ . The sum (resp. product) of  $\varphi$  and  $\varphi'$  is defined as the germ at the right of the origin of the sum (resp. product) of the semi-algebraic and continuous function  $f + f'$  (resp.  $ff'$ ) defined on  $(0, t)$ . It is easy to check that equipped with this addition and multiplication, the germs of semi-algebraic continuous functions at the right of the origin form a ring. The 0 (resp. 1) element of this ring is the germ of semi-algebraic continuous function at the right of the origin with representative the constant function with value 0 (resp. 1).

Consider a germ  $\varphi$  of semi-algebraic continuous function at the right of the origin and a representative  $f$  of  $\varphi$  defined on  $(0, t)$ . The set  $A = \{x \in (0, t) \mid f(x) = 0\}$  is a semi-algebraic set, and thus a finite union of points and intervals (Corollary 2.79). If  $A$  contains an interval  $(0, t')$ , then  $\varphi = 0$ . Otherwise, denoting by  $t'$  the smallest element of  $A$  (defined as  $t$  if  $A$  is empty), the restriction of  $f$  to  $(0, t')$  is everywhere non-zero, and hence  $1/f$  is a semi-algebraic and continuous function defined on  $(0, t')$  with associated germ  $1/\varphi$ . Thus the germs of semi-algebraic continuous functions at the right of the origin form a field.

Consider a germ  $\varphi$  of semi-algebraic continuous function at the right of the origin and a representative  $f$  of  $\varphi$  defined on  $(0, t)$ . The sets

$$\begin{aligned} A &= \{x \in (0, t) \mid f(x) = 0\}, \\ B &= \{x \in (0, t) \mid f(x) > 0\}, \\ C &= \{x \in (0, t) \mid f(x) < 0\}. \end{aligned}$$

are semi-algebraic and partition  $(0, t)$  into a finite number of points and intervals. One and only one of these three sets contains an interval of the form  $(0, t')$ . Thus, the sign of a germ  $\varphi$  of a semi-algebraic continuous function at the right of the origin is well defined. It is easy to check that equipped with this sign function, the germs of semi-algebraic continuous functions at the right of the origin form an ordered field.

It remains to prove that the germs of semi-algebraic continuous functions at the right of the origin have the intermediate value property, by Theorem 2.11.

It is sufficient to prove the intermediate value property for  $\bar{P}$  separable, by Lemma 3.12.

**Lemma 3.12.** *The property  $(I(P, a, b))$*

$$P(a)P(b) < 0 \Rightarrow \exists x \quad a < x < b \quad P(x) = 0$$

*holds for any  $P \in \mathbb{R}[X]$  if and only if it holds for any  $P \in \mathbb{R}[X]$ , with  $P$  separable.*

**Proof of Lemma 3.12:** It is clear that if  $(I(P, a, b))$  holds for any  $P \in \mathbb{R}[X]$ , it holds for any  $P \in \mathbb{R}[X]$ , with  $P$  separable. In the other direction, if  $P$  is separable, there is nothing to prove. So, suppose that  $P(a)P(b) < 0$ . If  $P_1 = \text{gcd}(P(X), P'(X)) \neq 1$ ,  $P(X) = P_1(X)P_2(X)$  with

$$\deg(P_1(X)) < \deg(P(X)), \deg(P_2(X)) < \deg(P(X)),$$

and either  $P_1(a)P_1(b) < 0$  or  $P_2(a)P_2(b) < 0$ . This process can be continued up to the moment where a divisor  $Q$  of  $P$ , with  $\text{gcd}(Q(X), Q'(X)) = 1$ ,  $Q(a)Q(b) < 0$  is found. Applying property  $(I(Q, a, b))$  gives a root of  $P$ .  $\square$



So, let  $\bar{P}(Y) = \alpha_p Y^p + \dots + \alpha_0$ ,  $\alpha_p \neq 0$ , be a separable polynomial, where the  $\alpha_i$  are germs of semi-algebraic continuous functions at the right of the origin, and let  $\varphi_1$  and  $\varphi_2$  be such that  $\bar{P}(\varphi_1) \bar{P}(\varphi_2) < 0$ . Let  $a_p, \dots, a_0, f_1, f_2$  be representatives of  $\alpha_p, \dots, \alpha_0, \varphi_1, \varphi_2$  defined on  $(0, t_0)$ . For every  $t \in (0, t_0)$ , let  $P(t, Y) = a_p(t) Y^p + \dots + a_0(t)$ . Shrinking  $(0, t_0)$ , if necessary, so that all the coefficients appearing in the signed remainder sequence of  $\bar{P}$ ,  $\bar{P}'$  have representatives defined on  $(0, t_0)$ , we can suppose that for every  $t \in (0, t_0)$ ,  $\deg(P(t, Y)) = p$ ,  $P(t, f_1(t))P(t, f_2(t)) < 0$ , and  $\gcd(P(t, Y), \bar{P}'(t, Y)) = 1$ . It is clear that, for every  $t \in (0, t_0)$ ,  $P(t, Y)$  has a root in  $(f_1(t), f_2(t))$ . Consider, for every  $0 \leq r \leq p$ , the set  $A_r \subset (0, t_0)$  of those  $t$  such that  $P(t, Y)$  has exactly  $r$  distinct roots in  $\mathbb{R}$ . Since  $A_r$  can be described by a formula, it is a semi-algebraic subset of  $(0, t_0)$ . The  $A_r$  partition  $(0, t_0)$  into a finite union of points and intervals, and exactly one of the  $A_r$  contains an interval of the form  $(0, t_1)$ . We are going to prove that for  $0 \leq i \leq r$ , the function  $g_i$  associating to  $t \in (0, t_1)$  the  $i$ -th root of  $P(t, Y)$  is semi-algebraic and continuous and that one of them lies between  $f_1$  and  $f_2$ .

Let  $t \in (0, t_1)$  and consider the  $g_i(t)$ . By Proposition 3.10, there exists an open interval  $(t - m, t + m)$  and semi-algebraic continuous functions  $h_i$  defined on  $(t - m, t + m)$  such that  $h_i(u)$  is a simple root of  $P(u, Y)$  for every  $u \in (t - m, t + m)$ . This root is necessarily  $g_i(u)$  because the number of roots of  $P(t, Y)$  on  $S$  is fixed. Thus,  $g_i$  is continuous.

Since for every  $t \in (0, t_1)$ ,  $P(t, f_1(t))P(t, f_2(t)) < 0$ , the graph of  $g_i$  does not intersect the graphs of  $f_1$  and  $f_2$ . So there is at least one  $g_i$  lying between  $f_1$  and  $f_2$ .  $\square$

**Proposition 3.13.** *The germs of semi-algebraic continuous functions at the right of the origin is the real closure of  $\mathbb{R}(\varepsilon)$  equipped with the unique order making  $\varepsilon$  infinitesimal. The element  $\varepsilon$  is sent to the germ of the identity map at the right of the origin.*

**Proof:** By Proposition 3.11, the germs of semi-algebraic continuous functions at the right of the origin form a real closed field. By Proposition 2.86, a germ of semi-algebraic function at the right of the origin is algebraic over  $\mathbb{R}(\varepsilon)$ .  $\square$

Using Corollary 2.98 and Proposition 3.13,

**Theorem 3.14.** *The real closed field of germs of semi-algebraic continuous functions at the right of the origin is isomorphic to the field of algebraic Puiseux series  $\mathbb{R}\langle\varepsilon\rangle$ .*

Using germs of semi-algebraic continuous functions at the right of the origin, the extension of a semi-algebraic set from  $\mathbb{R}$  to  $\mathbb{R}\langle\varepsilon\rangle$  has a particularly simple meaning. Before explaining this, we need a notation.

**Notation 3.15. [Composition with germs]** Consider a germ  $\varphi$  of semi-algebraic continuous functions at the right of the origin and  $f$  defined on  $(0, t)$  representing  $\varphi$ . If  $g$  is a continuous semi-algebraic function defined on the image of  $f$ , we denote by  $g \circ \varphi$  the germ of semi-algebraic continuous functions at the right of the origin associated to the semi-algebraic continuous function  $g \circ f$  defined on  $(0, t)$ . Note that  $g \circ \varphi$  is independent of the choice of the representative  $f$  of  $\varphi$ . Note also that if  $f$  represents  $\varphi$ ,  $f \circ \varepsilon = \varphi$ , since  $\varepsilon$  is the germ of the identity map at the right of the origin. □

**Proposition 3.16.** *Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set and  $\varphi = (\varphi_1, \dots, \varphi_k) \in \text{R}\langle\varepsilon\rangle^k$ . Let  $f_1, \dots, f_k$  be continuous semi-algebraic functions defined on  $(0, t)$  and representing  $\varphi_1, \dots, \varphi_k$  and let  $f = (f_1, \dots, f_k)$ . Then*

$$\varphi \in \text{Ext}(S, \text{R}\langle\varepsilon\rangle) \Leftrightarrow \exists t > 0 \quad \forall t' \quad 0 < t' < t \quad f(t') \in S.$$

*Suppose that  $\varphi \in \text{Ext}(S, \text{R}\langle\varepsilon\rangle)$  and let  $g$  be a semi-algebraic function defined on  $S$ . Then  $\text{Ext}(g, \text{R}\langle\varepsilon\rangle)(\varphi) = g \circ \varphi$ .*

*In particular,  $\text{Ext}(f, \text{R}\langle\varepsilon\rangle)(\varepsilon) = \varphi$ .*

**Proof:** The first part of the proposition is clear since, as we have seen above in the proof of Proposition 3.11, if  $P \in \text{R}[X_1, \dots, X_k]$  and  $\varphi_1, \dots, \varphi_k$  are germs of semi-algebraic continuous functions at the right of the origin with representatives  $f_1, \dots, f_k$  defined on a common  $(0, t)$ ,

- $P(\varphi_1, \dots, \varphi_k) = 0$  in  $\text{R}\langle\varepsilon\rangle$  if and only if there is an interval  $(0, t) \subset \mathbb{R}$  such that  $\forall t' \in (0, t) \quad P(f_1(t'), \dots, f_k(t')) = 0$
- $P(\varphi_1, \dots, \varphi_k) > 0$  in  $\text{R}\langle\varepsilon\rangle$  if and only if there is an interval  $(0, t) \subset \mathbb{R}$  such that  $\forall t' \in (0, t) \quad P(f_1(t'), \dots, f_k(t')) > 0$ .

The second part is clear as well by definition of the extension. The last part is a consequence of the second one, taking  $S = \text{R}\langle\varepsilon\rangle$ ,  $\varphi = \varepsilon$ ,  $f = \text{Id}$ ,  $g = f$  and using the remark at the end of Notation 3.15. □

An important property of  $\text{R}\langle\varepsilon\rangle$  is that sentences with coefficients in  $\text{R}[\varepsilon]$  which are true in  $\text{R}\langle\varepsilon\rangle$  are also true on a sufficiently small interval  $(0, r) \subset \mathbb{R}$ . Namely:

**Proposition 3.17.** *If  $\Phi$  is a sentence in the language of ordered fields with coefficients in  $\text{R}[\varepsilon]$  and  $\Phi'(t)$  is the sentence obtained by substituting  $t \in \mathbb{R}$  for  $\varepsilon$  in  $\Phi$ , then  $\Phi$  is true in  $\text{R}\langle\varepsilon\rangle$  if and only if there exists  $t_0$  in  $\mathbb{R}$  such that  $\Phi'(t)$  is true for every  $t \in (0, t_0) \cap \mathbb{R}$ .*

**Proof:** The semi-algebraic set  $A = \{t \in \mathbb{R} \mid \Phi'(t)\}$  is a finite union of points and intervals. If  $A$  contains an interval  $(0, t_0)$  with  $t_0$  a positive element of  $\mathbb{R}$ , then the extension of  $A$  to  $\text{R}\langle\varepsilon\rangle$  contains  $(0, t_0) \subset \text{R}\langle\varepsilon\rangle$ , so that  $\varepsilon \in \text{Ext}(A, \text{R}\langle\varepsilon\rangle)$  and  $\Phi = \Phi'(\varepsilon)$  is true in  $\text{R}\langle\varepsilon\rangle$ .

On the other hand, if  $A$  contains no interval  $(0, t)$  with  $t$  a positive element of  $\mathbb{R}$ , there exists  $t_0$  such that  $(0, t_0) \cap A = \emptyset$  and thus  $\text{Ext}((0, t_0) \cap A, \mathbb{R}\langle\varepsilon\rangle) = \emptyset$  and  $\varepsilon \notin \text{Ext}(A, \mathbb{R}\langle\varepsilon\rangle)$ , which means that  $\Phi$  is not true in  $\mathbb{R}\langle\varepsilon\rangle$ .  $\square$

The subring of germs of semi-algebraic continuous functions at the right of the origin which are bounded by an element of  $\mathbb{R}$  coincides with the valuation ring  $\mathbb{R}\langle\varepsilon\rangle_b$  defined in Chapter 2 (Notation 2.100). Indeed, it is clear by Proposition 3.17 that a germ  $\varphi$  of semi-algebraic continuous functions at the right of the origin is **bounded** by an element of  $\mathbb{R}$  if and only if  $\varphi$  has a representative  $f$  defined on  $(0, t)$  which is bounded. Note that this property is independent of the representative  $f$  chosen for  $\varphi$ .

The ring homomorphism  $\lim_\varepsilon$  defined on  $\mathbb{R}\langle\varepsilon\rangle_b$  in Notation 2.100 has a useful consequence for semi-algebraic functions.

**Proposition 3.18.** *Let  $f: (0, a) \rightarrow \mathbb{R}$  be a continuous bounded semi-algebraic function. Then  $f$  can be continuously extended to a function  $\bar{f}$  on  $[0, a)$ .*

**Proof:** Let  $M$  bound the absolute value of  $f$  on  $(0, a)$ . Thus  $M$  bounds the germ of semi-algebraic continuous function  $\varphi \in \mathbb{R}\langle\varepsilon\rangle$  associated to  $f$  using Proposition 3.16 and  $\lim_\varepsilon(\varphi)$  is well-defined. Let  $b = \lim_\varepsilon(\varphi)$ . Defining

$$\bar{f}(t) = \begin{cases} b & \text{if } t = 0, \\ f(t) & \text{if } t \in (0, a) \end{cases}$$

we easily see that  $\bar{f}$  is continuous at 0. Indeed for every  $r > 0$  in  $\mathbb{R}$ , the extension of the set  $\{t \in \mathbb{R} \mid |f(t) - b| \leq r\}$  to  $\mathbb{R}\langle\varepsilon\rangle$  contains  $\varepsilon$ , since  $\text{Ext}(f, \mathbb{R}\langle\varepsilon\rangle)(\varepsilon) - b = \varphi - b$  is infinitesimal, and therefore there is a positive  $\delta$  in  $\mathbb{R}$  such that it contains the interval  $(0, \delta)$  by Proposition 3.17.  $\square$

We can now prove a more geometric result. Note that its statement does not involve Puiseux series, while the proof we present does.

**Theorem 3.19. [Curve selection lemma]** *Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set. Let  $x \in \bar{S}$ . Then there exists a continuous semi-algebraic mapping  $\gamma: [0, 1) \rightarrow \mathbb{R}^k$  such that  $\gamma(0) = x$  and  $\gamma((0, 1)) \subset S$ .*

**Proof:** Let  $x \in \bar{S}$ . For every  $r > 0$  in  $\mathbb{R}$ ,  $B(x, r) \cap S$  is non-empty, hence  $B(x, \varepsilon) \cap \text{Ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  is non-empty by the Transfer principle (Theorem 2.80). Let  $\varphi \in B(x, \varepsilon) \cap \text{Ext}(S, \mathbb{R}\langle\varepsilon\rangle)$ . By Proposition 3.16 there exists a representative of  $\varphi$  which is a semi-algebraic continuous function  $f$  defined on  $(0, t)$  such that for every  $t', 0 < t' < t$ ,  $f(t') \in B(x, r) \cap S$ . Using Proposition 3.18 and scaling, we get  $\gamma: [0, 1) \rightarrow \mathbb{R}^k$  such that  $\gamma(0) = x$  and  $\gamma((0, 1)) \subset S$ . It is easy to check that  $\gamma$  is continuous at 0.  $\square$

### 3.4 Closed and Bounded Semi-algebraic Sets

In  $\mathbb{R}^k$ , a closed bounded set  $S$  is compact, i.e. has the property that whenever  $S$  is covered by a family of sets open in  $S$ , it is also covered by a finite subfamily of these sets. This is no longer true for a general real closed field  $\mathbb{R}$ , as can be seen by the following examples.

a) The interval  $[0, 1] \subset \mathbb{R}_{\text{alg}}$  is not compact since the family

$$\{[0, r) \cup (s, 1] \mid 0 < r < \pi/4 < s < 1, r \in \mathbb{R}_{\text{alg}}\}$$

(where  $\pi = 3.14\dots$ ), is an open cover of  $[0, 1]$  which has no finite subcover.

b) The interval  $[0, 1] \subset \mathbb{R}_{\text{alg}}$  is not compact since the family

$$\{[0, r) \cup (s, 1] \mid 0 < r < \pi/4 < s < 1, r \in \mathbb{R}_{\text{alg}}\}$$

(where  $\pi = 3.14\dots$ ), is an open cover of  $[0, 1]$  which has no finite subcover.

c) The interval  $[0, 1] \subset \mathbb{R}(\varepsilon)$  is not compact since the family

$$\{[0, f) \cup (r, 1] \mid f > 0 \text{ and infinitesimal over } \mathbb{R}, r \in \mathbb{R}, 0 < r < 1\}$$

is an open cover with no finite subcover.

However, closed and bounded semi-algebraic sets do enjoy properties of compact subsets, as we see now. We are going to prove the following result.

**Theorem 3.20.** *Let  $S$  be a closed, bounded semi-algebraic set and  $g$  a semi-algebraic continuous function defined on  $S$ . Then  $g(S)$  is closed and bounded.*

Though the statement of this theorem is geometric, the proof we present uses the properties of the real closed extension  $\mathbb{R}(\varepsilon)$  of  $\mathbb{R}$ .

The proof of the theorem uses the following lemma:

**Lemma 3.21.** *Let  $g$  be a semi-algebraic continuous function defined on a closed, bounded semi-algebraic set  $S \subset \mathbb{R}^k$ . If  $\varphi \in \text{Ext}(S, \mathbb{R}(\varepsilon))$ , then  $g \circ \varphi$  is bounded over  $\mathbb{R}$  and*

$$g(\lim_{\varepsilon}(\varphi)) = \lim_{\varepsilon}(g \circ \varphi).$$

**Proof:** Let  $f = (f_1, \dots, f_k)$  be a semi-algebraic function defined on  $(0, t)$  and representing  $\varphi = (\varphi_1, \dots, \varphi_k) \in \mathbb{R}(\varepsilon)^k$  and let  $\bar{f}$  its extension to  $[0, t)$ , using Proposition 3.18. By definition of  $\lim_{\varepsilon}$ ,

$$\bar{f}(0) = b = \lim_{\varepsilon}(\varphi)$$

since  $\varphi - b$  is infinitesimal. Since  $S$  is closed  $b \in S$ . Thus  $g$  is continuous at  $b$ . Hence, for every  $r > 0 \in \mathbb{R}$ , there is an  $\eta$  such that if  $z \in S$  and  $\|z - b\| < \eta$  then  $\|g(z) - g(b)\| < r$ . Using the Transfer Principle (Theorem 2.80) together with the fact that  $\varphi \in \text{Ext}(S, \mathbb{R}(\varepsilon))$  and  $\varphi - b$  is infinitesimal over  $\mathbb{R}$  we see that  $\|g \circ \varphi - g(b)\|$  is smaller than any  $r > 0$ . Thus  $g \circ \varphi$  is bounded over  $\mathbb{R}$  and infinitesimally close to  $g(b)$ , and hence  $g(\lim_{\varepsilon}(\varphi)) = \lim_{\varepsilon}(g \circ \varphi)$ . □

**Proof of Theorem 3.20:** We first prove that  $g(S)$  is closed. Suppose that  $x$  is in the closure of  $g(S)$ . Then  $B(x, r) \cap g(S)$  is not empty, for any  $r \in \mathbb{R}$ . Hence, by the Transfer principle (Theorem 2.80),  $B(x, \varepsilon) \cap \text{Ext}(g(S), \mathbb{R}\langle\varepsilon\rangle)$  is not empty. Thus, there is a  $\varphi \in \text{Ext}(g(S), \mathbb{R}\langle\varepsilon\rangle)$  for which  $\lim_\varepsilon(\varphi) = x$ . By Proposition 2.90, there is a  $\varphi' \in \text{Ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  such that  $g \circ \varphi' = \varphi$ . Since  $S$  is closed and bounded and  $\varphi'$  has a representative  $f'$  defined on  $(0, t)$  which can be extended continuously to  $\overline{f'}$  at 0,  $\lim_\varepsilon(\varphi') = \overline{f'}(0) \in S$ , and we conclude that  $g(\lim_\varepsilon(\varphi')) = \lim_\varepsilon(\varphi) = x$ . Hence  $x \in g(S)$ .

We now prove that  $g(S)$  is bounded. The set

$$A = \{t \in \mathbb{R} \mid \exists x \in S \quad \|g(x)\| = t\}$$

is semi-algebraic and so it is a finite union of points and intervals. For every  $\varphi \in \text{Ext}(S, \mathbb{R}\langle\varepsilon\rangle)$ ,  $g \circ \varphi$  is bounded over  $\mathbb{R}$  by Lemma 3.21. Thus  $\text{Ext}(A, \mathbb{R}\langle\varepsilon\rangle)$  does not contain  $1/\varepsilon$ . This implies that  $A$  contains no interval of the form  $(M, +\infty)$ , and thus  $A$  is bounded.  $\square$

### 3.5 Implicit Function Theorem

The usual notions of differentiability over  $\mathbb{R}$  can be developed over an arbitrary real closed field  $\mathbb{R}$ . We do this now.

Let  $f$  be a semi-algebraic function from a semi-algebraic open subset  $U$  of  $\mathbb{R}^k$  to  $\mathbb{R}^p$ , and let  $x_0 \in U$ . We write  $\lim_{x \rightarrow x_0} f(x) = y_0$  for

$$\forall r > 0 \exists \delta \forall x \quad \|x - x_0\| < \delta \Rightarrow \|f(x) - y_0\| < r$$

and  $f(x) = o(\|x - x_0\|)$  for

$$\lim_{x \rightarrow x_0} \frac{f(x)}{\|x - x_0\|} = 0.$$

If  $M$  is a semi-algebraic subset of  $U$ , we write  $\lim_{x \in M, x \rightarrow x_0} f(x) = y_0$  for

$$\forall r > 0 \exists \delta \forall x \in M \quad \|x - x_0\| < \delta \Rightarrow \|f(x) - y_0\| < r.$$

The function  $f: (a, b) \rightarrow \mathbb{R}$  is **differentiable at**  $x_0 \in (a, b)$  with derivative  $f'(x_0)$  if

$$\lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = f'(x_0).$$

We consider only semi-algebraic functions. Theorem 3.20 implies that a semi-algebraic function continuous on a closed and bounded interval is bounded and attains its bounds.

**Exercise 3.4.** Prove that Rolle's Theorem and the Mean Value Theorem hold for semi-algebraic differentiable functions.

**Proposition 3.22.** *Let  $f: (a, b) \rightarrow \mathbb{R}$  be a semi-algebraic function differentiable on the interval  $(a, b)$ . Then its derivative  $f'$  is a semi-algebraic function.*

**Proof:** Describe the graph of  $f'$  by a formula in the language of ordered fields with parameters in  $\mathbb{R}$ , and use Corollary 2.78. □

**Exercise 3.5.** Provide the details of the proof of Proposition 3.22.

Partial derivatives of multivariate semi-algebraic functions are defined in the usual way and have the usual properties. In particular let  $U \subset \mathbb{R}^k$  be a semi-algebraic open set and  $f: U \rightarrow \mathbb{R}^p$ , and suppose that the partial derivatives of the coordinate functions of  $f$  with respect to  $X_1, \dots, X_k$  exist on  $U$  and are continuous. These partial derivatives are clearly semi-algebraic functions.

For every  $x_0 \in U$ , let  $df(x_0)$  denote the **derivative of  $f$  at  $x_0$** , i.e. the linear mapping from  $\mathbb{R}^k$  to  $\mathbb{R}^p$  sending  $(h_1, \dots, h_k)$  to

$$\left( \sum_{j=1, \dots, k} \frac{\partial f_1}{\partial X_j}(x_0) h_j, \dots, \sum_{j=1, \dots, k} \frac{\partial f_p}{\partial X_j}(x_0) h_j \right).$$

The matrix of  $df(x_0)$  is the **Jacobian matrix of  $f$  at  $x_0$**  and its determinant is the **Jacobian of  $f$  at  $x_0$** . Following the usual arguments from a calculus course, It is clear that

$$f(x) - f(x_0) - df(x_0)(x - x_0) = o(\|x - x_0\|).$$

As in the univariate case, one can iterate the above definition to define higher derivatives.

Let  $U \subset \mathbb{R}^k$  be a semi-algebraic open set and  $B \subset \mathbb{R}^p$  a semi-algebraic set. The set of semi-algebraic functions from  $U$  to  $B$  for which all partial derivatives up to order  $\ell$  exist and are continuous is denoted  $\mathcal{S}^\ell(U, B)$ , and the class  $\mathcal{S}^\infty(U, B)$  is the intersection of  $\mathcal{S}^\ell(U, B)$  for all finite  $\ell$ . The ring  $\mathcal{S}^\ell(U, \mathbb{R})$  is abbreviated  $\mathcal{S}^\ell(U)$ , and the ring  $\mathcal{S}^\infty(U, \mathbb{R})$  is also called the ring of **Nash functions**.

We present a semi-algebraic version of the implicit function theorem whose proof is essentially the same as the classical proofs.

Given a linear mapping  $F: \mathbb{R}^k \rightarrow \mathbb{R}^p$ , we define the **norm of  $F$**  by  $\|F\| = \sup(\{\|F(x)\| \mid \|x\| = 1\})$ . This is a well-defined element of  $\mathbb{R}$  by Theorem 3.20, since  $x \mapsto \|F(x)\|$  is a continuous semi-algebraic function and  $\{x \mid \|x\| = 1\}$  is a closed and bounded semi-algebraic set.

**Proposition 3.23.** *Let  $x$  and  $y$  be two points of  $\mathbb{R}^k$ ,  $U$  an open semi-algebraic set containing the segment  $[x, y]$ , and  $f \in \mathcal{S}^1(U, \mathbb{R}^\ell)$ . Then*

$$\|f(x) - f(y)\| \leq M \|x - y\|,$$

where  $M = \sup(\{\|df(z)\| \mid z \in [x, y]\})$  ( $M$  is well defined, by Theorem 3.20).

**Proof:** Define  $g(t) = f((1-t)x + ty)$  for  $t \in [0, 1]$ . Then  $\|g'(t)\| \leq M \|x - y\|$  for  $t \in [0, 1]$ . For any positive  $c \in \mathbb{R}$ , we define

$$A_c = \{t \in [0, 1] \mid \|g(t) - g(0)\| \leq M \|x - y\| t + ct\}$$

which is a closed semi-algebraic subset of  $[0, 1]$  containing 0. Let  $t_0$  be the largest element in  $A_c$ . Suppose  $t_0 \neq 1$ . We have

$$\|g(t_0) - g(0)\| \leq M\|x - y\|t_0 + ct_0.$$

Since  $\|g'(t_0)\| \leq M\|x - y\|$ , we can find  $r > 0$  in  $\mathbb{R}$  such that if  $t_0 < t < t_0 + r$ ,

$$\|g(t) - g(t_0)\| \leq M\|x - y\|(t - t_0) + c(t - t_0).$$

So, for  $t_0 < t < t_0 + r$ , by summing the two displayed inequalities, we have

$$\|g(t) - g(0)\| \leq M\|x - y\|t + ct,$$

which contradicts the maximality of  $t_0$ . Thus  $1 \in A_c$  for every  $c$ , which gives the result.  $\square$

**Proposition 3.24. [Inverse Function Theorem]** *Let  $U'$  be a semi-algebraic open neighborhood of the origin 0 of  $\mathbb{R}^k$ ,  $f \in \mathcal{S}^\ell(U', \mathbb{R}^k)$ ,  $\ell \geq 1$ , such that  $f(0) = 0$  and that  $df(0): \mathbb{R}^k \rightarrow \mathbb{R}^k$  is invertible. Then there exist semi-algebraic open neighborhoods  $U, V$  of 0 in  $\mathbb{R}^k$ ,  $U \subset U'$ , such that  $f|_U$  is a homeomorphism onto  $V$  and  $(f|_U)^{-1} \in \mathcal{S}^\ell(V, U)$ .*

**Proof:** We can suppose that  $df(0)$  is the identity  $\text{Id}$  of  $\mathbb{R}^k$  (by composing with  $df(0)^{-1}$ ). Take  $g = f - \text{Id}$ . Then  $dg(0) = 0$ , and there is  $r_1 \in \mathbb{R}$  such that  $\|dg(x)\| \leq \frac{1}{2}$  if  $x \in B_k(0, r_1)$ . By Proposition 3.23, if  $x, y \in B_k(0, r_1)$ , then:

$$\|f(x) - f(y) - (x - y)\| \leq \frac{1}{2}\|x - y\|$$

and thus

$$\frac{1}{2}\|x - y\| \leq \|f(x) - f(y)\| \leq \frac{3}{2}\|x - y\|,$$

using the triangle inequalities. This implies that  $f$  is injective on  $B_k(0, r_1)$ . We can find  $r_2 < r_1$  with  $df(x)$  invertible for  $x \in B_k(0, r_2)$ . Now we prove that  $f(B_k(0, r_2)) \supset B_k(0, r_2/4)$ . For  $y^0$  with  $\|y^0\| < r_2/4$ , define  $h(x) = \|f(x) - y^0\|^2$ . Then  $h$  reaches its minimum on  $\overline{B_k(0, r_2)}$  and does not reach it on the boundary  $S^{k-1}(0, r_2)$  since if  $\|x\| = r_2$ , one has  $\|f(x)\| \geq r_2/2$  and thus  $h(x) > (r_2/4)^2 > h(0)$ . Therefore, this minimum is reached at a point  $x^0 \in B_k(0, r_2)$ . One then has, for  $i = 1, \dots, n$ ,

$$\frac{\partial h}{\partial x_i}(x^0) = 0, \quad \text{i.e.} \quad \sum_{j=1}^k (f_j(x^0) - y_j^0) \frac{\partial f_j}{\partial x_i}(x^0) = 0.$$

Since  $df(x^0)$  is invertible, we have  $f(x^0) = y^0$ . We then define  $V = B_k(0, r_2/4)$ ,  $U = f^{-1}(V) \cap B_k(0, r_2)$ . The function  $f^{-1}$  is continuous because

$$\|f^{-1}(x) - f^{-1}(y)\| \leq 2\|x - y\|$$

for  $x, y \in V$ , and we easily get  $d(f^{-1})(x) = (df(f^{-1}(x)))^{-1}$ .  $\square$

**Theorem 3.25. [Implicit Function Theorem]** *Let  $(x^0, y^0) \in \mathbb{R}^{k+\ell}$ , and let  $f_1, \dots, f_\ell$  be semi-algebraic functions of class  $\mathcal{S}^m$  on an open neighborhood of  $(x^0, y^0)$  such that  $f_j(x^0, y^0) = 0$  for  $j = 1, \dots, \ell$  and the Jacobian matrix of  $f = (f_1, \dots, f_\ell)$  at  $(x^0, y^0)$  with respect to the variables  $y_1, \dots, y_\ell$  is invertible. Then there exists a semi-algebraic open neighborhood  $U$  (resp.  $V$ ) of  $x^0$  (resp.  $y^0$ ) in  $\mathbb{R}^k$  (resp.  $\mathbb{R}^\ell$ ) and a function  $\varphi \in \mathcal{S}^m(U, V)$  such that  $\varphi(x^0) = y^0$ , and, for every  $(x, y) \in U \times V$ , we have*

$$f_1(x, y) = \dots = f_\ell(x, y) = 0 \Leftrightarrow y = \varphi(x).$$

**Proof:** Apply Proposition 3.24 to the function  $(x, y) \mapsto (x, f(x, y))$ . □

We now have all the tools needed to develop “semi-algebraic differential geometry”.

The notion of an  $\mathcal{S}^\infty$ -diffeomorphism between semi-algebraic open sets of  $\mathbb{R}^k$  is clear. The semi-algebraic version of  $\mathcal{C}^\infty$  submanifolds of  $\mathbb{R}^k$  is as follows.

An  $\mathcal{S}^\infty$ -**diffeomorphism**  $\varphi$  from a semi-algebraic open  $U$  of  $\mathbb{R}^k$  to a semi-algebraic open  $\Omega$  of  $\mathbb{R}^k$  is a bijection from  $U$  to  $\Omega$  that is  $\mathcal{S}^\infty$  and such that  $\varphi^{(-1)}$  is  $\mathcal{S}^\infty$ .

A semi-algebraic subset  $M$  of  $\mathbb{R}^k$  is an  $\mathcal{S}^\infty$  **submanifold of  $\mathbb{R}^k$  of dimension  $\ell$**  if for every point  $x$  of  $M$ , there exists a semi-algebraic open  $U$  of  $\mathbb{R}^k$  and an  $\mathcal{S}^\infty$ -diffeomorphism  $\varphi$  from  $U$  to a semi-algebraic open neighborhood  $\Omega$  of  $x$  in  $\mathbb{R}^k$  such that  $\varphi(0) = x$  and

$$\varphi(U \cap (\mathbb{R}^\ell \times \{0\})) = M \cap \Omega$$

(where  $\mathbb{R}^\ell \times \{0\} = \{(a_1, \dots, a_\ell, 0, \dots, 0) \mid (a_1, \dots, a_\ell) \in \mathbb{R}^\ell\}$ ).

A semi-algebraic map from  $M$  to  $N$ , where  $M$  (resp.  $N$ ) is an  $\mathcal{S}^\infty$  submanifold of  $\mathbb{R}^m$  (resp.  $\mathbb{R}^n$ ), is an  $\mathcal{S}^\infty$  **map** if it is locally the restriction of an  $\mathcal{S}^\infty$  map from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ .

A point  $x$  of a semi-algebraic set  $S \subset \mathbb{R}^k$  is a **smooth point of dimension  $\ell$**  if there is a semi-algebraic open subset  $U$  of  $S$  containing  $x$  which is an  $\mathcal{S}^\infty$  submanifold of  $\mathbb{R}^k$  of dimension  $\ell$ .

Let  $x$  be a smooth point of dimension  $\ell$  of an  $\mathcal{S}^\infty$  submanifold  $M$  of  $\mathbb{R}^k$  and let  $\Omega$  be a semi-algebraic open neighborhood of  $x$  in  $\mathbb{R}^k$  and  $\varphi: U \rightarrow \Omega$  as in the definition of a submanifold. Let  $X_1, \dots, X_k$  be the coordinates of the domain of  $\varphi = (\varphi_1, \dots, \varphi_k)$ . We call the set  $T_x(M) = x + d\varphi(0)(\mathbb{R}^\ell \times \{0\})$  the **tangent space to  $M$  at  $x$** . Clearly, the tangent space contains  $x$  and is a translate of an  $\ell$  dimensional linear subspace of  $\mathbb{R}^k$ , i.e. an  **$\ell$ -flat**. More concretely, note that the tangent space  $T_x(M)$  is the translate by  $x$  of the linear space spanned by the first  $\ell$  columns of the Jacobian matrix.

We next prove the usual geometric properties of tangent spaces.

**Proposition 3.26.** *Let  $x$  be a point of an  $\mathcal{S}^\infty$  submanifold  $M$  of  $\mathbb{R}^k$  having dimension  $\ell$  and let  $\pi$  denote orthogonal projection onto the  $\ell$ -flat  $T_x(M)$ .*

*Then,  $\lim_{y \in M, y \rightarrow x} \frac{\|y - \pi(y)\|}{\|y - x\|} = 0$ .*



**Proof:** Let  $\Omega$  be a semi-algebraic open neighborhood of  $x$  in  $\mathbb{R}^k$  and  $\varphi: U \rightarrow \Omega$  as in the definition of a submanifold. Let  $X_1, \dots, X_k$  be the coordinates of the domain of  $\varphi = (\varphi_1, \dots, \varphi_k)$ . Then,

$$T_x(M) = x + d\varphi(0)(\mathbb{R}^\ell \times \{0\}).$$

From elementary properties of derivatives (see Equation (3.5)), it is clear that for  $u \in \mathbb{R}^\ell \times \{0\}$ ,  $\varphi(u) - d\varphi(0)(u) = o(\|u\|)$ .

Now, for  $y \in M \cap \Omega$ , let  $u = \varphi^{-1}(y)$ . Then, since  $\pi$  is an orthogonal projection,

$$\|y - \pi(y)\| \leq \|\varphi(u) - d\varphi(0)(u)\| = o(\|u\|).$$

Since,  $\varphi^{-1}$  is an  $\mathcal{S}^\infty$  map, for any bounded neighborhood of  $x$  there is a constant  $C$  such that  $\|\varphi^{-1}(y)\| \leq C\|y - x\|$  for all  $y$  in the neighborhood. Since  $\|u\| = \|\varphi^{-1}(y)\| \leq C\|y - x\|$ ,

$$\|\varphi(u) - d\varphi(0)(u)\| = o(\|y - x\|),$$

and the conclusion follows.  $\square$

We next prove that the tangent vector at a point of a curve lying on an  $\mathcal{S}^\infty$  submanifold  $M$  of  $\mathbb{R}^k$  is contained in the tangent space to  $M$  at that point.

**Proposition 3.27.** *Let  $x$  be a point of the  $\mathcal{S}^\infty$  submanifold  $M$  in  $\mathbb{R}^k$  having dimension  $\ell$ , and let  $\gamma: [-1, 1] \rightarrow \mathbb{R}^k$  be an  $\mathcal{S}^\infty$  curve contained in  $M$  with  $\gamma(0) = x$ . Then the tangent vector  $x + \gamma'(0)$  is contained in the tangent space  $T_x(M)$ .*

**Proof:** Let  $\gamma(t) = (\gamma_1(t), \dots, \gamma_k(t))$ . Let  $\Omega, \varphi$  be as in the definition of submanifold, and consider the composite map  $\varphi^{-1} \circ \gamma: [-1, 1] \rightarrow \mathbb{R}^k$ . Applying the chain rule,  $d(\varphi^{-1} \circ \gamma)(0) = d\varphi^{-1}(x)(\gamma'(0))$ . Since  $\gamma([-1, 1]) \subset M$ , it follows that  $\varphi^{-1}(\gamma([-1, 1])) \subset \mathbb{R}^\ell \times \{0\}$ , and  $d(\varphi^{-1} \circ \gamma)(t) \in \mathbb{R}^\ell \times \{0\}$  for all  $t \in [-1, 1]$ . Thus,  $d\varphi^{-1}(x)(\gamma'(0)) \in \mathbb{R}^\ell \times \{0\}$ . Since  $d\varphi^{-1}(x) = (d\varphi(0))^{-1}$ , applying  $d\varphi(0)$  to both sides we have  $\gamma'(0) \in d\varphi(0)(\mathbb{R}^\ell \times \{0\})$ , and finally  $x + \gamma'(0) \in T_x(M)$ .  $\square$

The notion of derivatives defined earlier for multivariate functions can now be extended to  $\mathcal{S}^\infty$  submanifolds.

Let  $f: M \rightarrow N$  be an  $\mathcal{S}^\infty$  map, where  $M$  (resp.  $N$ ) is a  $m'$  (resp.  $n'$ ) dimensional  $\mathcal{S}^\infty$  submanifold of  $\mathbb{R}^m$  (resp.  $\mathbb{R}^n$ ).

Let  $x \in M$  and let  $\Omega$  (resp.  $\Omega'$ ) be a neighborhood of  $x$  (resp.  $f(x)$ ) in  $\mathbb{R}^m$  (resp.  $\mathbb{R}^n$ ) and  $\varphi$  (resp.  $\psi$ ) a semi-algebraic diffeomorphism from  $U$  to  $\Omega$  (resp.  $U'$  to  $\Omega'$ ) such that  $\varphi(0) = x$  (resp.  $\psi(0) = f(x)$ ) and

$$\varphi(\mathbb{R}^{m'} \times \{0\}) = M \cap \Omega \quad (\text{resp. } \psi(\mathbb{R}^{n'} \times \{0\}) = N \cap \Omega').$$

Clearly,  $\psi^{-1} \circ f \circ \varphi: \mathbb{R}^m \rightarrow \mathbb{R}^n$  is an  $\mathcal{S}^\infty$  map, and its restriction to  $\mathbb{R}^{m'} \times \{0\}$  is an  $\mathcal{S}^\infty$  map to  $\mathbb{R}^{n'} \times \{0\}$ .

The derivative  $d(\psi^{-1} \circ f \circ \varphi)(0)$  restricted to  $\mathbb{R}^{m'} \times \{0\}$  maps  $\mathbb{R}^{m'} \times \{0\}$  into  $\mathbb{R}^{n'} \times \{0\}$ .

The linear map  $df(x): T_x(M) \rightarrow T_{f(x)}(N)$  defined by

$$df(x)(v) = f(x) + d\psi(0)(d(\psi^{-1} \circ f \circ \varphi)(0)(d\varphi^{-1}(x)(v - x))),$$

is called the **derivative** of  $f$  at  $x$ .

**Proposition 3.28.**

- a) *A semi-algebraic open subset of an  $\mathcal{S}^\infty$  submanifold  $V$  of dimension  $i$  is an  $\mathcal{S}^\infty$  submanifold of dimension  $i$ .*
- b) *If  $V'$  is an  $\mathcal{S}^\infty$  submanifold of dimension  $j$  contained in an  $\mathcal{S}^\infty$  submanifold  $V$  of dimension  $i$ , then  $j \leq i$ .*

**Proof:** a) is clear. b) follows from the fact that the tangent space to  $V'$  at  $x \in V'$  is a subspace of the tangent space to  $V$  at  $x$ . □

### 3.6 Bibliographical Notes

Semi-algebraic sets appear first in a logical context in Tarski's work [154]. They were studied from a geometrical and topological point of view by Brakhage [28], in his unpublished thesis. The modern study of semi-algebraic sets starts with Lojasiewicz, as a particular case of semi-analytic sets [110, 111].

## Algebra

---

We start in Section 4.1 with the discriminant, and the related notion of subdiscriminant. In Section 4.2, we define the resultant and signed subresultant coefficients of two univariate polynomials and indicate how to use them for real root counting. We describe in Section 4.3 an algebraic real root counting technique based on the signature of a quadratic form. We then give a constructive proof of Hilbert's Nullstellensatz using resultants in Section 4.4. In Section 4.5, we algebraically characterize systems of polynomials with a finite number of solutions and prove that the corresponding quotient rings are finite dimensional vector spaces. In Section 4.6, we give a multivariate generalization of the real root counting technique based on the signature of a quadratic form described in Section 4.3. In Section 4.7, we define projective space and prove a weak version of Bézout's theorem.

Throughout Chapter 4,  $K$  is a field of characteristic zero and  $C$  is an algebraically closed field containing it. We will also denote by  $R$  a real closed field containing  $K$  when  $K$  is an ordered field.

### 4.1 Discriminant and Subdiscriminant

**Notation 4.1. [Discriminant]** Let  $P \in R[X]$  be a monic polynomial of degree  $p$ ,

$$P = X^p + a_{p-1}X^{p-1} + \dots + a_0,$$

and let  $x_1, \dots, x_p$  be the roots of  $P$  in  $C$  (repeated according to their multiplicities). The **discriminant** of  $P$ ,  $\text{Disc}(P)$ , is defined by

$$\text{Disc}(P) = \prod_{p \geq i > j \geq 1} (x_i - x_j)^2. \quad \square$$

*Remark 4.2.* The discriminant played a key role in the algebraic proof of the fundamental theorem of algebra (proof of  $a \Rightarrow b$ ) in Theorem 2.11, see Remark 2.17). □

**Proposition 4.3.**  $\text{Disc}(P) = 0$  if and only if  $\deg(\text{gcd}(P, P')) > 0$ .

**Proof:** It is clear from the definition that  $\text{Disc}(P) = 0$  if and only if  $P$  has a multiple root in  $\mathbb{C}$ .  $\square$

*Remark 4.4.* When all the roots of  $P$  are in  $\mathbb{R}$  and distinct,  $\text{Disc}(P) > 0$ .  $\square$

The sign of the discriminant counts the number of real roots modulo 4.

**Proposition 4.5.** *Let  $P \in \mathbb{R}[X]$  be monic with  $\mathbb{R}$  real closed, of degree  $p$ , and with  $p$  distinct roots in  $\mathbb{C}$ ; Denoting by  $t$  the number of roots of  $P$  in  $\mathbb{R}$ ,*

$$\begin{aligned}\text{Disc}(P) > 0 &\Leftrightarrow t \equiv p \pmod{4}, \\ \text{Disc}(P) < 0 &\Leftrightarrow t \equiv p - 2 \pmod{4}.\end{aligned}$$

**Proof:** Let  $y_1, \dots, y_t$  be the roots of  $P$  in  $\mathbb{R}$  and  $z_1, \bar{z}_1, \dots, z_s, \bar{z}_s$  the roots of  $P$  in  $\mathbb{C} \setminus \mathbb{R}$ , with  $\mathbb{C} = \mathbb{R}[i]$ .

The conclusion is clear since

$$\begin{aligned}\text{sign}\left(\prod_{i=1}^s (z_i - \bar{z}_i)^2\right) &= (-1)^s, \\ (y_i - y_j)^2 &> 0, 1 \leq i < j \leq t, \\ ((z_i - z_j)(z_i - \bar{z}_j)(\bar{z}_i - z_j)(\bar{z}_i - \bar{z}_j))^2 &> 0, 1 \leq i < j \leq s, \\ ((y_i - z_j)(y_i - \bar{z}_j))^2 &> 0, 1 \leq i \leq t, 1 \leq j \leq s.\end{aligned}$$

Thus,  $\text{Disc}(P) > 0$  if and only if  $s$  is even, and  $\text{Disc}(P) < 0$  if and only if  $s$  is odd.  $\square$

The  $p - k$ -**subdiscriminant** of  $P$ ,  $1 \leq k \leq p$ , is by definition

$$\text{sDisc}_{p-k}(P) = \sum_{\substack{I \subset \{1, \dots, p\} \\ \#(I) = k}} \prod_{\substack{(j, \ell) \in I \\ \ell > j}} (x_j - x_\ell)^2.$$

Note that  $\text{sDisc}_{p-1}(P) = p$ . The discriminant is the 0-th subdiscriminant:

$$\text{sDisc}_0(P) = \text{Disc}(P) = \prod_{p \geq j > \ell \geq 1} (x_j - x_\ell)^2.$$

*Remark 4.6.* It is clear that when all the roots of  $P$  are in  $\mathbb{R}$

$$\text{sDisc}_0(P) = \dots = \text{sDisc}_{j-1}(P) = 0, \text{sDisc}_j(P) \neq 0$$

if and only if  $P$  has  $p - j$  distinct roots. We shall see later in Proposition 4.29 that this property is true in general.  $\square$

The subdiscriminants are intimately related to the Newton sums of  $P$ .

**Definition 4.7.** The  $i$ -**th Newton sum** of the polynomial  $P$ , denoted  $N_i$ , is  $\sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x)x^i$ , where  $\mu(x)$  is the multiplicity of  $x$ .  $\square$

The Newton sums can be obtained from the coefficients of  $P$  by the famous Newton identities.

**Proposition 4.8.** *Let  $P = a_p X^p + a_{p-1} X^{p-1} + \dots + a_1 X + a_0$ . For any  $i$*

$$(p - i) a_{p-i} = a_p N_i + \dots + a_0 N_{i-p}, \tag{4.1}$$

*with the convention  $a_i = N_i = 0$  for  $i < 0$ .*

**Proof:** We have

$$P = a_p \prod_{x \in \text{Zer}(P, \mathbb{C})} (X - x)^{\mu(x)},$$

$$\frac{P'}{P} = \sum_{x \in \text{Zer}(P, \mathbb{C})} \frac{\mu(x)}{(X - x)}.$$

Using

$$\frac{1}{X - x} = \sum_{i=0}^{\infty} \frac{x^i}{X^{i+1}},$$

we get

$$\frac{P'}{P} = \sum_{i=0}^{\infty} \frac{N_i}{X^{i+1}},$$

$$P' = \left( \sum_{i=0}^{\infty} \frac{N_i}{X^{i+1}} \right) P.$$

Equation (4.1) follows by equating the coefficients of  $X^{p-i-1}$  on both sides of the last equality. □

Consider the square matrix

$$\text{Newt}_{p-k}(P) = \begin{bmatrix} N_0 & N_1 & \dots & & \dots & N_{k-1} \\ N_1 & \dots & & \dots & N_{k-1} & N_k \\ \dots & & \dots & N_{k-1} & N_k & \dots \\ & \dots & N_{k-1} & N_k & \dots & \\ \dots & N_{k-1} & N_k & \dots & & \dots \\ N_{k-1} & N_k & \dots & & \dots & N_{2k-2} \end{bmatrix}$$

with entries the Newton sums of the monic polynomial  $P$  of degree  $p$ .

We denote as usual by  $\det(M)$  the determinant of a square matrix  $M$ .

**Proposition 4.9.** *For every  $k$ ,  $1 \leq k \leq p$ ,*

$$\text{sDisc}_{p-k}(P) = \det(\text{Newt}_{p-k}(P)).$$

The proof of Proposition 4.9 uses the Cauchy-Binet formula.

**Proposition 4.10. [Cauchy-Binet]** Let  $A$  be a  $n \times m$  matrix and  $B$  be a  $m \times n$  matrix,  $m \geq n$ . For every  $I \subset \{1, \dots, m\}$  of cardinality  $n$ , denote by  $A_I$  the  $n \times n$  matrix obtained by extracting from  $A$  the columns with indices in  $I$ . Similarly let  $B^I$  be the  $n \times n$  matrix obtained by extracting from  $B$  the rows with indices in  $I$ .

$$\det(AB) = \sum_{\substack{I \subset \{1, \dots, m\} \\ \#(I) = n}} \det(A_I) \det(B^I).$$

**Proof:**

We introduce an  $m$ -dimensional diagonal matrix  $D_\lambda$  with diagonal entries the variables  $\lambda_1, \dots, \lambda_m$  and study  $\det(AD_\lambda B)$ . Since the entries of the matrix  $AD_\lambda B$  are homogeneous linear forms in the  $\lambda_i$ ,  $\det(AD_\lambda B)$  is a homogeneous polynomial of degree  $n$  in the  $\lambda_i$ .

We are going to prove that the only monomials with non-zero coefficients of  $\det(AD_\lambda B)$  are of the form  $\lambda_I = \prod_{i \in I} \lambda_i$  for a subset  $I \subset \{1, \dots, m\}$ ,  $\#(I) = n$ .

Indeed if we consider  $I \subset \{1, \dots, m\}$ ,  $\#(I) < n$ , the specialization of  $\det(AD_\lambda B)$  obtained by sending  $\lambda_j$  to 0 for  $j \notin I$  is identically null. This implies that the coefficients of all the monomials where a variable is repeated are 0.

If we choose  $I \subset \{1, \dots, m\}$ ,  $\#(I) = n$ , and specialize the variables  $\lambda_i$ ,  $i \in I$  to 1 and the variables  $\lambda_i$ ,  $i \notin I$  to 0, we get the coefficient of  $\lambda_I = \prod_{i \in I} \lambda_i$  in  $\det(AD_\lambda B)$ , which is  $\det(A_I) \det(B^I)$ .

Specializing finally all the  $\lambda_i$  to 1, we get the required identity.  $\square$

The proof of Proposition 4.9 makes also use of the classical Vandermonde determinant. Let  $x_1, \dots, x_r$  be elements of a field  $K$ . The **Vandermonde determinant** of  $x_1, \dots, x_r$  is  $\det(V(x_1, \dots, x_r))$  with

$$V(x_1, \dots, x_{r-1}, x_r) = \begin{bmatrix} 1 & \cdots & 1 & 1 \\ x_1 & \cdots & x_{r-1} & x_r \\ \vdots & & \vdots & \vdots \\ x_1^{r-1} & \cdots & x_{r-1}^{r-1} & x_r^{r-1} \end{bmatrix}$$

the **Vandermonde matrix**.

**Lemma 4.11.**

$$\det(V(x_1, \dots, x_r)) = \prod_{r \geq i > j \geq 1} (x_i - x_j).$$

**Proof:** The claim is true when  $x_1, \dots, x_r$  are not all distinct since both sides are 0. The proof when  $x_1, \dots, x_r$  are all distinct is by induction on  $r$ . The claim is obviously true for  $r = 2$ . Suppose that the claim is true for  $r - 1$  and consider

$$V(x_1, \dots, x_{r-1}, X) = \begin{bmatrix} 1 & \cdots & 1 & 1 \\ x_1 & \cdots & x_{r-1} & X \\ \vdots & & \vdots & \vdots \\ x_1^{r-1} & \cdots & x_{r-1}^{r-1} & X^{r-1} \end{bmatrix}.$$

The polynomial  $\det(V(x_1, \dots, x_{r-1}, X))$  has degree at most  $r - 1$ , with  $r - 1$  distinct roots  $x_1, \dots, x_{r-1}$  because, replacing  $X$  by  $x_i$  in  $V(x_1, \dots, x_{r-1}, X)$ , we get a matrix with two equal columns. So

$$\det(V(x_1, \dots, x_{r-1}, X)) = c \prod_{r-1 \geq j \geq 1} (X - x_j).$$

The coefficient of  $\det(V(x_1, \dots, x_{r-1}, X))$  is the Vandermonde determinant of  $x_1, \dots, x_{r-1}$ ,  $\det(V(x_1, \dots, x_{r-1}))$  is equal to

$$\prod_{r-1 \geq i > j \geq 1} (x_i - x_j),$$

by the induction hypothesis. So

$$\det(V(x_1, \dots, x_{r-1}, X)) = \prod_{r-1 \geq i > j \geq 1} (x_i - x_j) \prod_{r-1 \geq j \geq 1} (X - x_j).$$

Now substitute  $x_r$  for  $X$  to get the claim. □

**Proof of Proposition 4.9:** Define

$$V_k = \begin{bmatrix} 1 & \dots & \dots & \dots & 1 \\ x_1 & \dots & \dots & \dots & x_p \\ \vdots & & & & \vdots \\ x_1^{k-1} & \dots & \dots & \dots & x_p^{k-1} \end{bmatrix}.$$

It is clear that  $V_k V_k^t = \text{Newt}_{p-k}(P)$ . Now apply Binet-Cauchy formula, noting that, if  $I \subset \{1, \dots, p\}$ ,  $\#(I) = k$ , and  $V_{kI}$  is the  $k \times k$  matrix obtained by extracting from  $V_k$  the columns with indices in  $I$

$$\det(V_{kI}) = \prod_{\substack{(j, \ell) \in I \\ \ell > j}} (x_j - x_\ell),$$

by Lemma 4.11. □

## 4.2 Resultant and Subresultant Coefficients

### 4.2.1 Resultant

Let  $P$  and  $Q$  be two non-zero polynomials of degree  $p$  and  $q$  in  $D[X]$ , where  $D$  is a ring. When  $D$  is a domain, its fraction field is denoted by  $K$ . Let

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + \dots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \dots + b_0. \end{aligned}$$

We define the Sylvester matrix associated to  $P$  and  $Q$  and the resultant of  $P$  and  $Q$ .

**Notation 4.12. [Sylvester matrix]** The **Sylvester matrix** of  $P$  and  $Q$ , denoted by  $\text{Syl}(P, Q)$ , is the matrix

$$\begin{bmatrix} a_p & \cdots & \cdots & \cdots & \cdots & a_0 & 0 & \cdots & 0 \\ 0 & \ddots & & & & & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & & & & & \ddots & 0 \\ 0 & \cdots & 0 & a_p & \cdots & \cdots & \cdots & \cdots & a_0 \\ b_q & \cdots & \cdots & \cdots & b_0 & 0 & \cdots & \cdots & 0 \\ 0 & \ddots & & & & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & & & & \ddots & \ddots & \vdots \\ \vdots & & & & & & & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & b_q & \cdots & \cdots & \cdots & b_0 \end{bmatrix}.$$

It has  $p + q$  columns and  $p + q$  rows. Note that its rows are

$$X^{q-1}P, \dots, P, X^{p-1}Q, \dots, Q$$

considered as vectors in the basis  $X^{p+q-1}, \dots, X, 1$ .

The **resultant** of  $P$  and  $Q$ , denoted  $\text{Res}(P, Q)$ , is the determinant of  $\text{Syl}(P, Q)$ .  $\square$

*Remark 4.13.* This matrix comes about quite naturally since it is the transpose of the matrix of the linear mapping  $U, V \mapsto UP + VQ$ , where  $(U, V)$  is identified with

$$(u_{q-1}, \dots, u_0, v_{p-1}, \dots, v_0),$$

and  $U = u_{q-1}X^{q-1} + \dots + u_0$ ,  $V = v_{p-1}X^{p-1} + \dots + v_0$ .  $\square$

The following lemma is clear from this remark.

**Lemma 4.14.** *Let  $D$  be a domain. Then  $\text{Res}(P, Q) = 0$  if and only if there exist non-zero polynomials  $U \in K[X]$  and  $V \in K[X]$ , with  $\deg(U) < q$  and  $\deg(V) < p$ , such that  $UP + VQ = 0$ .*

We can now prove the well-known proposition.

**Proposition 4.15.** *Let  $D$  be a domain. Then  $\text{Res}(P, Q) = 0$  if and only if  $P$  and  $Q$  have a common factor in  $K[X]$ .*

**Proof:** The proposition is an immediate consequence of the preceding lemma and of Proposition 1.5, since the least common multiple of  $P$  and  $Q$  has degree  $< p + q$  if and only if there exist non-zero polynomials  $U$  and  $V$  with  $\deg(U) < q$  and  $\deg(V) < p$  such that  $UP + VQ = 0$ .  $\square$



If  $D$  is a domain, with fraction field  $K$ ,  $a_p \neq 0$  and  $b_q \neq 0$ , the resultant can be expressed as a function of the roots of  $P$  and  $Q$  in an algebraically closed field  $C$  containing  $K$ .

**Theorem 4.16.** *Let*

$$P = a_p \prod_{i=1}^p (X - x_i)$$

$$Q = b_q \prod_{j=1}^q (X - y_j),$$

*in other words  $x_1, \dots, x_p$  are the roots of  $P$  (counted with multiplicities) and  $y_1, \dots, y_q$  are the roots of  $Q$  (counted with multiplicities).*

$$\text{Res}(P, Q) = a_p^q b_q^p \prod_{i=1}^p \prod_{j=1}^q (x_i - y_j).$$

**Proof:** Let

$$\Theta(P, Q) = a_p^q b_q^p \prod_{i=1}^p \prod_{j=1}^q (x_i - y_j).$$

If  $P$  and  $Q$  have a root in common,  $\text{Res}(P, Q) = \Theta(P, Q) = 0$ , and the theorem holds. So we suppose now that  $P$  and  $Q$  are coprime. The theorem is proved by induction on the length  $n$  of the remainder sequence of  $P$  and  $Q$ .

When  $n = 2$ ,  $Q$  is a constant  $b$ , and  $\text{Res}(P, Q) = \Theta(P, Q) = b^p$ .

The induction step is based on the following lemma.

**Lemma 4.17.** *Let  $R$  be the remainder of the Euclidean division of  $P$  by  $Q$  and let  $r$  be the degree of  $R$ . Then,*

$$\text{Res}(P, Q) = (-1)^{pq} b_q^{p-r} \text{Res}(Q, R),$$

$$\Theta(P, Q) = (-1)^{pq} b_q^{p-r} \Theta(Q, R).$$

**Proof of Lemma 4.17:** Let  $R = c_r X^r + \dots + c_0$ . Replacing the rows of coefficients of the polynomials  $X^{q-1}P, \dots, P$  by the rows of coefficients of the polynomials  $X^{q-1}R, \dots, R$  in the Sylvester matrix of  $P$  and  $Q$  gives the matrix

$$M = \begin{bmatrix} 0 & 0 & c_r & \cdots & \cdots & c_0 & 0 & \cdots & 0 \\ \vdots & & \ddots & \ddots & & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & & & \ddots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 & c_r & \cdots & \cdots & c_0 \\ b_q & \cdots & \cdots & \cdots & b_0 & 0 & \cdots & \cdots & 0 \\ 0 & \ddots & & & & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & & & & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & & & & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & b_q & \cdots & \cdots & \cdots & b_0 \end{bmatrix}.$$

such that

$$\det(M) = \text{Res}(P, Q).$$

Indeed,

$$R = P - \sum_{i=0}^{p-q} d_i (X^i Q),$$

where  $C = \sum_{i=0}^{p-q} d_i X^i$  is the quotient of  $P$  in the euclidean division of  $P$  by  $Q$ , and adding to a row a multiple of other rows does not change the determinant.

Denoting by  $N$  the matrix whose rows are  $X^{p-1} Q, \dots, X^{r-1} Q, \dots, Q, X^{q-1} R, \dots, R$ , we note that

$$N = \begin{bmatrix} b_q & \dots & \dots & \dots & b_0 & 0 & \dots & \dots & 0 \\ 0 & b_q & & & & b_0 & \ddots & & \vdots \\ \vdots & 0 & b_q & \dots & \dots & \dots & b_0 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & & & & \ddots & 0 \\ \vdots & \vdots & \dots & 0 & b_q & \dots & \dots & \dots & b_0 \\ \vdots & \vdots & c_r & \dots & \dots & c_0 & 0 & \dots & 0 \\ \vdots & \vdots & 0 & \ddots & & & \ddots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \ddots & & & \ddots & 0 \\ 0 & 0 & 0 & \dots & 0 & c_r & \dots & \dots & c_0 \end{bmatrix}.$$

is obtained from  $M$  by exchanging the order of rows, so that

$$\det(N) = (-1)^{pq} \det(M).$$

It is clear, developing the determinant of  $N$  by its  $p - r$  first columns, that

$$\det(N) = b_q^{p-r} \text{Res}(Q, R).$$

On the other hand, since  $P = CQ + R$ ,  $P(y_j) = R(y_j)$  and

$$\Theta(P, Q) = a_p^q \prod_{i=1}^p Q(x_i) = (-1)^{pq} b_q^p \prod_{j=1}^q P(y_j),$$

we have

$$\begin{aligned} \Theta(P, Q) &= (-1)^{pq} b_q^p \prod_{j=1}^q P(y_j) \\ &= (-1)^{pq} b_q^p \prod_{j=1}^q R(y_j) \\ &= (-1)^{pq} b_q^{p-r} \Theta(Q, R). \end{aligned}$$

□ □

For any ring  $D$ , the following holds:

**Proposition 4.18.** *If  $P, Q \in D[X]$ , then there exist  $U, V \in D[X]$  such that  $\deg(U) < q$ ,  $\deg(V) < p$ , and  $\text{Res}(P, Q) = UP + VQ$ .*

**Proof:** Let  $\text{Syl}(P, Q)^*$  be the matrix whose first  $p + q - 1$  columns are the first  $p + q - 1$  first columns of  $\text{Syl}(P, Q)$  and such that the elements of the last column are the polynomials  $X^{q-1}P, \dots, P, X^{p-1}Q, \dots, Q$ . Using the linearity of  $\det(\text{Syl}(P, Q)^*)$  as a function of its last column it is clear that

$$\det(\text{Syl}(P, Q)^*) = \text{Res}(P, Q) + \sum_{j=1}^{p+q-1} d_j X^j,$$

where  $d_j$  is the determinant of the matrix  $\text{Syl}(P, Q)_j$  whose first  $p + q - 1$  columns are the first  $p + q - 1$  columns of  $\text{Syl}(P, Q)$  and such that the last column is the  $p + q - j$ -th column of  $\text{Syl}(P, Q)$ . Since  $\text{Syl}(P, Q)_j$  has two identical columns,  $d_j = 0$  for  $j = 1, \dots, p + q - 1$  and

$$\det(\text{Syl}(P, Q)^*) = \text{Res}(P, Q).$$

Expanding the determinant of  $\text{Syl}(P, Q)^*$  by its last column, we obtain the claimed identity.  $\square$

The Sylvester matrix and the resultant also have the following useful interpretation. Let  $\mathbb{C}$  be an algebraically closed field. Identify a monic polynomial

$$X^q + b_{q-1}X^{q-1} + \dots + b_0 \in \mathbb{C}[X]$$

of degree  $q$  with the point  $(b_{q-1}, \dots, b_0) \in \mathbb{C}^q$ . Let

$$\begin{aligned} m: \mathbb{C}^q \times \mathbb{C}^p &\longrightarrow \mathbb{C}^{q+p} \\ (Q, P) &\longmapsto QP \end{aligned}$$

be the mapping defined by the multiplication of monic polynomials. The map  $m$  sends

$$(b_{q-1}, \dots, b_0, a_{p-1}, \dots, a_0)$$

to the vector whose entries are  $(m_{p+q-1}, \dots, m_0)$ , where

$$m_j = \sum_{q-i+p-k=j} b_{q-i} a_{p-k} \text{ for } j = p + q - 1, \dots, 0$$

(with  $b_q = a_p = 1$ ). The following proposition is thus clear:

**Proposition 4.19.** *The Jacobian matrix of  $m$  is the Sylvester matrix of  $P$  and  $Q$  and the Jacobian of  $m$  is the resultant.*

Finally, the definition of resultants as determinants implies that:

**Proposition 4.20.** *If  $P$  is monic,  $\deg(Q) \leq \deg(P)$ , and  $f: \mathbb{D} \rightarrow \mathbb{D}'$  is a ring homomorphism, then  $f(\text{Res}(P, Q)) = \text{Res}(f(P), f(Q))$  (denoting by  $f$  the induced homomorphism from  $\mathbb{D}[X]$  to  $\mathbb{D}'[X]$ ).*

### 4.2.2 Subresultant Coefficients

We now define the Sylvester-Habicht matrices and the signed subresultant coefficients of  $P$  and  $Q$  when  $p = \deg(P) > q = \deg(Q)$ .

**Notation 4.21. [Sylvester-Habicht matrix]** Let  $0 \leq j \leq q$ . The  $j$ -th **Sylvester-Habicht matrix** of  $P$  and  $Q$ , denoted  $\text{SyHa}_j(P, Q)$ , is the matrix whose rows are  $X^{q-j-1}P, \dots, P, Q, \dots, X^{p-j-1}Q$  considered as vectors in the basis  $X^{p+q-j-1}, \dots, X, 1$ :

$$\begin{bmatrix} a_p & \cdots & \cdots & \cdots & \cdots & a_0 & 0 & 0 \\ 0 & \ddots & & & & & \ddots & 0 \\ \vdots & \ddots & a_p & \cdots & \cdots & \cdots & \cdots & a_0 \\ \vdots & & 0 & b_q & \cdots & \cdots & \cdots & b_0 \\ \vdots & \ddots & \ddots & & & & \ddots & 0 \\ 0 & \ddots & & & & \ddots & \ddots & \vdots \\ b_q & \cdots & \cdots & \cdots & b_0 & 0 & \cdots & 0 \end{bmatrix}.$$

It has  $p + q - j$  columns and  $p + q - 2j$  rows.

The  $j$ -th **signed subresultant coefficient** denoted  $\text{sRes}_j(P, Q)$  or  $\text{sRes}_j$  is the determinant of the square matrix  $\text{SyHa}_{j,j}(P, Q)$  obtained by taking the first  $p + q - 2j$  columns of  $\text{SyHa}_j(P, Q)$ .

By convention, we extend these definitions for  $q < j \leq p$  by

$$\begin{aligned} \text{sRes}_p(P, Q) &= \text{sign}(a_p), \\ \text{sRes}_j(P, Q) &= 0, \quad q < j < p. \end{aligned}$$

□

*Remark 4.22.* The matrix  $\text{SyHa}_j(P, Q)$  comes about quite naturally since it is the transpose of the matrix of the mapping  $U, V \mapsto UP + VQ$ , where  $(U, V)$  is identified with

$$(u_{q-j-1}, \dots, u_0, v_0, \dots, v_{p-j-1}),$$

with  $U = u_{q-j-1}X^{q-j-1} + \dots + u_0$ ,  $V = v_{p-j-1}X^{p-j-1} + \dots + v_0$ .

The peculiar order of rows is adapted to the real root counting results presented later, in Chapter 8. □

The following lemma is clear from this remark:

**Lemma 4.23.** *Let  $D$  be a domain and  $0 \leq j \leq \min(p, q)$  if  $p \neq q$  (resp.  $0 \leq j \leq p - 1$  if  $p = q$ ). Then  $\text{sRes}_j(P, Q) = 0$  if and only if there exist non-zero polynomials  $U \in K[X]$  and  $V \in K[X]$ , with  $\deg(U) < q - j$  and  $\deg(V) < p - j$ , such that  $\deg(UP + VQ) < j$ .*

The following proposition will be useful for the cylindrical decomposition in Chapter 5.

**Proposition 4.24.** *Let  $D$  be a domain and  $0 \leq j \leq \min(p, q)$  if  $p \neq q$  (resp.  $0 \leq j \leq p-1$  if  $p = q$ ). Then  $\deg(\gcd(P, Q)) \geq j$  if and only if*

$$\text{sRes}_0(P, Q) = \cdots = \text{sRes}_{j-1}(P, Q) = 0.$$

**Proof:** Suppose that  $\deg(\gcd(P, Q)) \geq j$ . Then, the least common multiple of  $P$  and  $Q$ ,

$$\text{lcm}(P, Q) = \frac{PQ}{\gcd(P, Q)}$$

(see Proposition 1.5) has degree  $\leq p + q - j$ . This is clearly equivalent to the existence of polynomials  $U$  and  $V$ , with  $\deg(U) \leq q - j$  and  $\deg(V) \leq p - j$ , such that  $UP = -VQ = \text{lcm}(P, Q)$ . Or, equivalently, that there exist polynomials  $U$  and  $V$  with  $\deg(U) \leq q - j$  and  $\deg(V) \leq p - j$  such that  $UP + VQ = 0$ . This implies that

$$\text{sRes}_0 = \cdots = \text{sRes}_{j-1} = 0$$

using Lemma 4.23.

The reverse implication is proved by induction on  $j$ . If  $j = 1$ ,  $\text{sRes}_0 = 0$  implies, using Lemma 4.23, that there exist  $U$  and  $V$  with  $\deg(U) < q$  and  $\deg(V) < p$  satisfying  $UP + VQ = 0$ . Hence  $\deg(\gcd(P, Q)) \geq 1$ . If

$$\text{sRes}_0(P, Q) = \cdots = \text{sRes}_{j-2}(P, Q) = 0,$$

the induction hypothesis implies that  $\deg(\gcd(P, Q)) \geq j - 1$ . If in addition  $\text{sRes}_{j-1} = 0$  then, by Lemma 4.23, there exist  $U$  and  $V$  with  $\deg(U) \leq q - j$  and  $\deg(V) \leq p - j$  such that  $\deg(UP + VQ) < j - 1$ . Since the greatest common divisor of  $P$  and  $Q$  divides  $UP + VQ$  and has degree  $\geq j - 1$ , we have  $UP + VQ = 0$ , which implies that  $\deg(\text{lcm}(P, Q)) \leq p + q - j$  and hence  $\deg(\gcd(P, Q)) \geq j$ .  $\square$

The following consequence is clear, using Lemma 4.23 and Proposition 4.24.

**Proposition 4.25.** *Let  $D$  be a domain and  $0 \leq j \leq \min(p, q)$  if  $p \neq q$  (resp.  $0 \leq j \leq p - 1$  if  $p = q$ ). Then  $\deg(\gcd(P, Q)) = j$  if and only if*

$$\text{sRes}_0(P, Q) = \cdots = \text{sRes}_{j-1}(P, Q) = 0, \text{sRes}_j(P, Q) \neq 0.$$

**Notation 4.26. [Reversing rows]** We denote by  $\varepsilon_i$  the signature of the permutation reversing the order of  $i$  consecutive rows in a matrix, i.e.  $\varepsilon_i = (-1)^{i(i-1)/2}$ . For every natural number  $i \geq 1$ ,

$$\varepsilon_{4i} = 1, \varepsilon_{4i-1} = -1, \varepsilon_{4i-2} = -1, \varepsilon_{4i-3} = 1. \quad (4.2)$$

In particular,  $\varepsilon_{i-2j} = (-1)^j \varepsilon_i$ .  $\square$

Thus, it is clear from the definitions that

$$\text{sRes}_0(P, Q) = \varepsilon_p \text{Res}(P, Q). \quad (4.3)$$

Note that, as a consequence Proposition 4.15 is a special case of Proposition 4.24.

Let us make the connection between subresultant coefficients and subdiscriminants.

We first define subdiscriminants of non-monic polynomials. Let

$$\begin{aligned}
 P &= a_p X^p + \dots + a_0, \\
 \text{sDisc}_{p-k}(P) &= a_p^{2k-2} \text{sDisc}_{p-k}(P/a_p) \\
 &= a_p^{2k-2} \sum_{\substack{I \subset \{1, \dots, p\} \\ \#(I)=k}} \prod_{(j, \ell) \in I, \ell > j} (x_j - x_\ell)^2
 \end{aligned}$$

**Proposition 4.27.**

$$a_p \text{sDisc}_{p-k}(P) = \text{sRes}_{p-k}(P, P'). \tag{4.4}$$

**Proof:** Indeed if

$$D_k = \begin{bmatrix}
 1 & 0 & \dots & 0 & 0 & \dots & \dots & \dots & 0 \\
 0 & 1 & \ddots & \vdots & \vdots & & & & \vdots \\
 \vdots & \ddots & \ddots & 0 & \vdots & & & & \vdots \\
 \vdots & & \ddots & 1 & 0 & \dots & \dots & \dots & 0 \\
 0 & \dots & \dots & 0 & N_0 & N_1 & \dots & \dots & N_{k-1} \\
 \vdots & & \ddots & \ddots & \vdots & & & & \vdots \\
 \vdots & & \ddots & \ddots & \ddots & & & & \vdots \\
 0 & N_0 & N_1 & \dots & \dots & N_{k-1} & \dots & \dots & N_{2k-3} \\
 N_0 & N_1 & \dots & \dots & N_{k-1} & \dots & \dots & \dots & N_{2k-2}
 \end{bmatrix},$$

and

$$D'_k = \begin{bmatrix}
 a_p & \dots & \dots & \dots & \dots & \dots & \dots & \dots & a_{p-2k+2} \\
 0 & a_p & \ddots & & & & & & \vdots \\
 \vdots & \ddots & \ddots & \ddots & & & & & \vdots \\
 \vdots & & 0 & a_p & \ddots & & & & \vdots \\
 \vdots & & & 0 & a_p & \ddots & & & \vdots \\
 \vdots & & & & \ddots & \ddots & \ddots & & \vdots \\
 \vdots & & & & & \ddots & \ddots & \ddots & \vdots \\
 \vdots & & & & & & a_p & & \vdots \\
 0 & \dots & \dots & \dots & \dots & \dots & \dots & 0 & a_p
 \end{bmatrix},$$

it is a easy to see that  $\text{SyHa}_{p-k, p-k}(P, P') = D_k \cdot D'_k$ , using the relations (4.1). Since  $\det(D'_k) = a_p^{2k-1}$ ,

$$\begin{aligned}
 \det(\text{SyHa}_{p-k, p-k}(P, P')) &= a_p^{2k-1} \text{sDisc}_{p-k}(P/a_p) \\
 &= a_p \text{sDisc}_{p-k}(P). \\
 \det(\text{SyHa}_{p-k, p-k}(P, P')) &= a_p^{2k-1} \text{sDisc}_{p-k}(P/a_p) \\
 &= a_p \text{sDisc}_{p-k}(P).
 \end{aligned}$$

On the other hand  $\det(D_k) = \text{sRes}_{p-k}(P, P')$ . The claim follows by Proposition 4.18. □

*Remark 4.28.* Note that if  $P \in D[X]$ , then  $\text{sDisc}_i(P) \in D$  for every  $i \leq p$ .  $\square$

**Proposition 4.29.** *Let  $D$  be a domain. Then  $\deg(\gcd(P, P')) = j$ ,  $0 \leq j < p$  if and only if*

$$\text{sDisc}_0(P) = \dots = \text{sDisc}_{j-1}(P) = 0, \text{sDisc}_j(P) \neq 0.$$

**Proof:** Follows immediately from Proposition 4.27 and Proposition 4.25  $\square$

### 4.2.3 Subresultant Coefficients and Cauchy Index

We indicate how to compute the Cauchy index by using only the signed subresultant coefficients. We need a definition:

**Notation 4.30. [Generalized Permanences minus Variations]**

Let  $s = s_p, \dots, s_0$  be a finite list of elements in an ordered field  $K$  such that  $s_p \neq 0$ . Let  $q < p$  such that  $s_{p-1} = \dots = s_{q+1} = 0$ , and  $s_q \neq 0$ , and  $s' = s_q, \dots, s_0$ . (if there exist no such  $q$ ,  $s'$  is the empty list). We define inductively

$$\text{PmV}(s) = \begin{cases} 0 & \text{if } s' = \emptyset, \\ \text{PmV}(s') + \varepsilon_{p-q} \text{sign}(s_p s_q) & \text{if } p - q \text{ is odd,} \\ \text{PmV}(s') & \text{if } p - q \text{ is even.} \end{cases}$$

where  $\varepsilon_{p-q} = (-1)^{(p-q)(p-q-1)/2}$ , using Notation 4.26.

Note that when all elements of  $s$  are non-zero,  $\text{PmV}(s)$  is the difference between the number of sign permanence and the number of sign variations in  $s_p, \dots, s_0$ . Note also that when  $s$  is the sequence of coefficients of polynomials  $\mathcal{P} = P_p, \dots, P_0$  with  $\deg(P_i) = i$ , then

$$\text{PmV}(s) = \text{Var}(\mathcal{P}; -\infty, +\infty)$$

(see Notation 2.32).  $\square$

Let  $P$  and  $Q$  be two polynomials with:

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + \dots + a_0 \\ Q &= b_{p-1} X^{p-1} + \dots + b_0, \end{aligned}$$

$\deg(P) = p, \deg(Q) = q \leq p - 1$ .

We denote by  $\text{sRes}(P, Q)$  the sequence of  $\text{sRes}_j(P, Q)$ ,  $j = p, \dots, 0$ .

**Theorem 4.31.**  $\text{PmV}(\text{sRes}(P, Q)) = \text{Ind}(Q/P)$ .

Before proving Theorem 4.31 let us list some of its consequences.

**Theorem 4.32.** *Let  $P$  and  $Q$  be polynomials in  $D[X]$  and  $R$  the remainder of  $P'Q$  and  $P$ . Then  $\text{PmV}(\text{sRes}(P, R)) = \text{TaQ}(Q, P)$ .*

**Proof:** Apply Theorem 4.31 and Proposition 2.57, since

$$\text{Ind}(P'Q/P) = \text{Ind}(R/P)$$

by Remark 2.55. □

**Theorem 4.33.** *Let  $P$  be a polynomial in  $D[X]$ . Then*

$$\text{PmV}(\text{sDisc}_{p-1}(P), \dots, \text{sDisc}_0(P))$$

*is the number of roots of  $P$  in  $R$ .*

**Proof:** Apply Theorem 4.31 and Proposition 4.27. □

The proof of Theorem 4.31 uses the following two lemmas.

**Lemma 4.34.**

$$\text{Ind}(Q/P) = \begin{cases} \text{Ind}(-R/Q) + \text{sign}(a_p b_q) & \text{if } p-q \text{ is odd,} \\ \text{Ind}(-R/Q) & \text{if } p-q \text{ is even.} \end{cases}$$

**Proof:** The claim is an immediate consequence of Lemma 2.60. □

**Lemma 4.35.**

$$\text{PmV}(\text{sRes}(P, Q)) = \begin{cases} \text{PmV}(\text{sRes}(Q, -R)) + \text{sign}(a_p b_q) & \text{if } p-q \text{ is odd,} \\ \text{PmV}(\text{sRes}(Q, -R)) & \text{if } p-q \text{ is even.} \end{cases}$$

The proof of Lemma 4.35 is based on the following proposition.

**Proposition 4.36.** *Let  $r$  be the degree of  $R = \text{Rem}(P, Q)$ .*

$$\text{sRes}_j(P, Q) = \varepsilon_{p-q} b_q^{p-r} \text{sRes}_j(Q, -R) \text{ if } j \leq r,$$

where  $\varepsilon_i = (-1)^{i(i-1)/2}$ .

Moreover,  $\text{sRes}_j(P, Q) = \text{sRes}_j(Q, -R) = 0$  if  $r < j < q$ .

**Proof:** Replacing the polynomials  $X^{q-j-1}P, \dots, P$  by the polynomials  $X^{q-j-1}R, \dots, R$  in  $\text{SyHa}_j(P, Q)$  does not modify the determinant based on the  $p+q-2j$  first columns. Indeed,

$$R = P - \sum_{i=0}^{p-q} c_i X^i Q,$$

where  $C = \sum_{i=0}^{p-q} c_i X^i$  is the quotient of  $P$  in the euclidean division of  $P$  by  $Q$ , and adding to a polynomial of a sequence a multiple of another polynomial of the sequence does not change the determinant based on the  $p+q-2j$  first columns.



Reversing the order of the polynomials multiplies the determinant based on the  $p + q - 2j$  first columns. by  $\varepsilon_{p+q-2j}$ . Replacing  $R$  by  $-R$  multiplies the determinant based on the  $p + q - 2j$  first columns by  $(-1)^{q-j}$ , and

$$(-1)^{q-j} \varepsilon_{p+q-2j} = \varepsilon_{p-q}$$

(see Notation 4.26). Denoting by  $D_j$  the determinant obtained by taking the  $p + q - 2j$  first columns of the matrix the rows corresponding to the coefficients of  $X^{p-j-1}Q, \dots, Q, -R, \dots, -X^{q-j-1}R$ ,

$$\text{sRes}_j(P, Q) = \varepsilon_{p-q} D_j.$$

If  $j \leq r$ , it is clear that

$$D_j = b_q^{p-r} \text{sRes}_j(Q, -R).$$

If  $r < j < q$ , it is clear that

$$D_j = \text{sRes}_j(P, Q) = \text{sRes}_j(Q, -R) = 0.$$

using the convention in Notation 4.20 and noting that the  $q - j$ -th row of the determinant  $D_j$  is null.  $\square$

**Proof of Lemma 4.35:** Using Proposition 4.36,

$$\text{PmV}(\text{sRes}_r(P, Q), \dots, \text{sRes}_0(P, Q)) = \text{PmV}(\text{sRes}_r(Q, -R), \dots, \text{sRes}_0(Q, -R)).$$

If  $q - r$  is even

$$\begin{aligned} & \text{PmV}(\text{sRes}_q(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_r(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_r(Q, -R), \dots, \text{sRes}_0(Q, -R)) \\ &= \text{PmV}(\text{sRes}_q(Q, -R), \dots, \text{sRes}_0(Q, -R)). \end{aligned}$$

If  $q - r$  is odd, since

$$\begin{aligned} \text{sRes}_q(P, Q) &= \varepsilon_{p-q} b_q^{p-q}, \\ \text{sRes}_q(Q, -R) &= \text{sign}(b_q), \\ \text{sRes}_r(P, Q) &= \varepsilon_{p-q} b_q^{p-r} \text{sRes}_r(Q, -R), \end{aligned}$$

denoting  $d_r = \text{sRes}_r(Q, -R)$ ,

$$\begin{aligned} & \text{PmV}(\text{sRes}_q(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_r(P, Q), \dots, \text{sRes}_0(P, Q)) + \varepsilon_{q-r} \text{sign}(b_q d_r) \\ &= \text{PmV}(\text{sRes}_q(Q, -R), \dots, \text{sRes}_0(Q, -R)). \end{aligned}$$

Thus in all cases

$$\begin{aligned} & \text{PmV}(\text{sRes}_q(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_q(Q, -R), \dots, \text{sRes}_0(Q, -R)). \end{aligned}$$

If  $p - q$  is even

$$\begin{aligned} & \text{PmV}(\text{sRes}_p(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_q(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_q(Q, -R), \dots, \text{sRes}_0(Q, -R)). \end{aligned}$$

If  $p - q$  is odd, since

$$\begin{aligned} \text{sRes}_p(P, Q) &= \text{sign}(a_p), \\ \text{sRes}_q(P, Q) &= \varepsilon_{p-q} b_q^{p-q}, \\ \text{PmV}(\text{sRes}_p(P, Q), \dots, \text{sRes}_0(P, Q)) \\ &= \text{PmV}(\text{sRes}_q(P, Q), \dots, \text{sRes}_0(P, Q)) + \text{sign}(a_p b_q) \\ &= \text{PmV}(\text{sRes}_q(Q, -R), \dots, \text{sRes}_0(Q, -R)) + \text{sign}(a_p b_q). \end{aligned}$$

□

**Proof of Theorem 4.31:** The proof proceeds by induction on the number  $n$  of elements with distinct degrees in the signed subresultant sequence.

If  $n = 2$ ,  $Q$  divides  $P$ . We have

$$\text{Ind}(Q/P) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p-q \text{ is odd,} \\ 0 & \text{if } p-q \text{ is even.} \end{cases}$$

by Lemma 4.34 and

$$\text{PmV}(\text{sRes}(P, Q)) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p-q \text{ is odd,} \\ 0 & \text{if } p-q \text{ is even.} \end{cases}$$

by Lemma 4.35.

Let us suppose that the theorem holds for  $n - 1$  and consider  $P$  and  $Q$  such that their signed subresultant sequence has  $n$  elements with distinct degrees. The signed subresultant sequence of  $Q$  and  $-R$  has  $n - 1$  elements with distinct degrees. By the induction hypothesis,

$$\text{PmV}(\text{sRes}(Q, -R)) = \text{Ind}(-R/Q).$$

So, by Lemma 4.34 and Lemma 4.35,

$$\text{PmV}(\text{sRes}(P, Q)) = \text{Ind}(Q/P). \quad \square$$

*Example 4.37.* Consider again  $P = X^4 + aX^2 + bX + c$ ,

$$\begin{aligned} \text{sDisc}_3(P) &= 4, \\ \text{sDisc}_2(P) &= -8a, \\ \text{sDisc}_1(P) &= 4(8ac - 9b^2 - 2a^3) \\ \text{sDisc}_0(P) &= 256c^3 - 128a^2c^2 + 144ab^2c + 16a^4c - 27b^4 - 4a^3b^2. \end{aligned}$$

As in Example 1.15, let

$$s = 8ac - 9b^2 - 2a^3,$$

$$\delta = 256c^3 - 128a^2c^2 + 144ab^2c + 16a^4c - 27b^4 - 4a^3b^2.$$

We indicate in the following tables the number of real roots of  $P$  (computed using Theorem 4.31) in the various cases corresponding to all the possible signs for  $a, s, \delta$ :

1	+	+	+	+	+	+	+	+	+
4	+	+	+	+	+	+	+	+	+
$-a$	+	+	+	+	+	+	+	+	+
$s$	+	+	+	-	-	-	0	0	0
$\delta$	+	-	0	+	-	0	+	-	0
$n$	4	2	3	0	2	1	2	2	2

1	+	+	+	+	+	+	+	+	+
4	+	+	+	+	+	+	+	+	+
$-a$	-	-	-	-	-	-	-	-	-
$s$	+	+	+	-	-	-	0	0	0
$\delta$	+	-	0	+	-	0	+	-	0
$n$	0	-2	-1	0	2	1	0	0	0

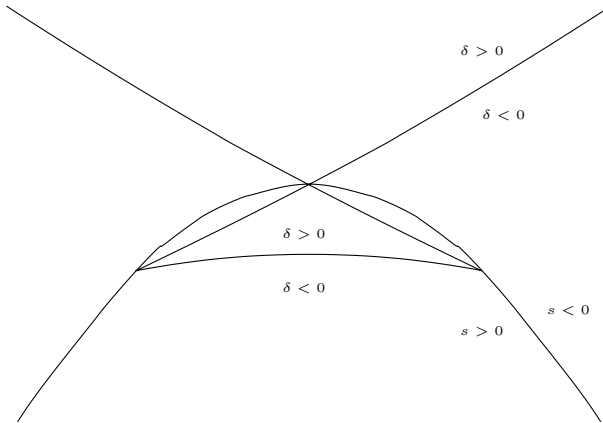
1	+	+	+	+	+	+	+	+	+
4	+	+	+	+	+	+	+	+	+
$-a$	0	0	0	0	0	0	0	0	0
$s$	+	+	+	-	-	-	0	0	0
$\delta$	+	-	0	+	-	0	+	-	0
$n$	2	0	1	0	2	1	0	2	1

Note that when  $a = s = 0$ , according to the definition of PmV when there are two consecutive zeroes,

$$\begin{cases} \text{PmV}(\text{sRes}(P, P')) = 0 & \text{if } \delta > 0 \\ \text{PmV}(\text{sRes}(P, P')) = 2 & \text{if } \delta < 0 \\ \text{PmV}(\text{sRes}(P, P')) = 1 & \text{if } \delta = 0. \end{cases}$$

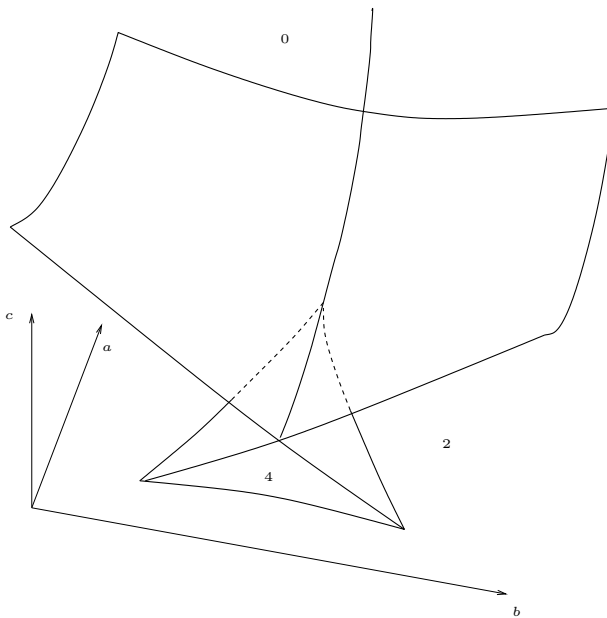
Notice that the only sign conditions on  $a, s, \delta$  for which all the roots of  $P$  are real is  $a < 0, s > 0, \delta > 0$ , according to Corollary 9.8. Remark that, according to Corollary 4.3, when  $\delta < 0$  there are always two distinct real roots. This looks incompatible with the tables we just gave. In fact, the sign conditions with  $\delta < 0$  giving a number of real roots different from 2, and the sign conditions with  $\delta > 0$  giving a number of real roots equal to 2 have empty realizations.

We represent in Figure 4.1 the set of polynomials of degree 4 in the plane  $a = -1$  and the zero sets of  $s, \delta$ .



**Fig. 4.1.**  $a = -1, s = \delta = 0$

Finally, in Figure 4.2 we represent the set of polynomials of degree 4 in  $a, b, c$  space and the zero sets of  $s, \delta$ .



**Fig. 4.2.** The set defined by  $\delta = 0$  and the different regions labelled by the number of real roots

**Exercise 4.1.** Find all sign conditions on  $a, s, \delta$  with non-empty realizations.

As a consequence, the formula  $\exists X \quad X^4 + aX^2 + bX + c = 0$  is equivalent to the quantifier-free formula

$$\begin{aligned} & (a < 0 \wedge s \geq 0 \wedge \delta > 0) \vee \\ & (a < 0 \wedge \delta \leq 0) \vee \\ & (a > 0 \wedge s < 0 \wedge \delta \leq 0) \vee \\ & (a = 0 \wedge s > 0 \wedge \delta \geq 0) \vee \\ & (a = 0 \wedge s \leq 0 \wedge \delta \leq 0). \end{aligned}$$

collecting all sign conditions giving  $n \geq 1$ . It can be checked easily that the realization of the sign conditions  $(a = 0 \wedge s > 0 \wedge \delta \geq 0)$  and  $(a < 0 \wedge s = 0 \wedge \delta > 0)$  are empty. So that  $(\exists X) X^4 + aX^2 + bX + c = 0$  is finally equivalent to

$$\begin{aligned} & (a < 0 \wedge s > 0 \wedge \delta > 0) \vee \\ & (a < 0 \wedge \delta \leq 0) \vee \\ & (a > 0 \wedge s < 0 \wedge \delta \leq 0) \vee \\ & (a = 0 \wedge s \leq 0 \wedge \delta \leq 0). \end{aligned}$$

It is interesting to compare this result with Example 2.63: the present description is more compact and involves only sign conditions on the principal subresultants  $a, s, \delta$ .  $\square$

## 4.3 Quadratic Forms and Root Counting

### 4.3.1 Quadratic Forms

The **transpose** of an  $n \times m$  matrix  $A = [a_{i,j}]$  is the  $m \times n$  matrix  $A^t = [b_{j,i}]$  defined by  $b_{j,i} = a_{i,j}$ . A square matrix  $A$  is **symmetric** if  $A^t = A$ .

A **quadratic form** with coefficients in a field  $K$  of characteristic 0 is a homogeneous polynomial of degree 2 in a finite number of variables of the form

$$\Phi(f_1, \dots, f_n) = \sum_{i,j=1}^n m_{i,j} f_i f_j$$

with  $M = [m_{i,j}]$  a symmetric matrix of size  $n$ . If  $f = (f_1, \dots, f_n)$ , then  $\Phi = f \cdot M \cdot f^t$ , where  $f^t$  is the transpose of  $f$ . The **rank** of  $\Phi$ , denoted by  $\text{Rank}(\Phi(f_1, \dots, f_n))$ , is the rank of the matrix  $M$ .

A **diagonal expression** of the quadratic form  $\Phi(f_1, \dots, f_n)$  is an identity

$$\Phi(f_1, \dots, f_n) = \sum_{i=1}^r c_i L_i(f_1, \dots, f_n)^2$$

with  $c_i \in K, c_i \neq 0$  and the  $L_i(f_1, \dots, f_n)$  are linearly independent linear forms with coefficients in  $K$ . The elements  $c_i, i = 1, \dots, r$  are the **coefficients** of the diagonal expression. Note that  $r = \text{Rank}(\Phi(f_1, \dots, f_n))$ .

**Theorem 4.38. [Sylvester's law of inertia]**

- A quadratic form  $\Phi(f_1, \dots, f_n)$  of dimension  $n$  has always a diagonal expression.
- If  $\mathbb{K}$  is ordered, the difference between the number of positive coefficients and the number of negative coefficients in a diagonal expression of  $\Phi(f_1, \dots, f_n)$  is a well defined quantity.

**Proof:** Let  $\Phi(f_1, \dots, f_n) = \sum_{i,j=1}^n m_{i,j} f_i f_j$ .

The first claim is proved by induction on  $n$ . The result is obviously true if  $n = 1$ . It is also true when  $M = 0$ .

If some diagonal entry  $m_{i,i}$  of  $M$  is not zero, we can suppose without loss of generality (reordering the variables) that  $m_{n,n}$  is not 0. Take

$$L(f_1, \dots, f_n) = \sum_{k=1}^n m_{k,n} f_k.$$

The quadratic form

$$\Phi(f_1, \dots, f_n) - \frac{1}{m_{n,n}} L(f_1, \dots, f_n)^2$$

does not depend on the variable  $f_n$ , and we can apply the induction hypothesis to

$$\Phi_1(f_1, \dots, f_{n-1}) = \Phi(f_1, \dots, f_n) - \frac{1}{m_{n,n}} L(f_1, \dots, f_n)^2.$$

Since  $L(f_1, \dots, f_n)$  is a linear form containing  $f_n$ , it is certainly linearly independent from the linear forms in the decomposition of  $\Phi_1(f_1, \dots, f_{n-1})$ .

If all diagonal entries are equal to 0, but  $M \neq 0$ , we can suppose without loss of generality (reordering the variables) that  $m_{n-1,n} \neq 0$ . Performing the linear change of variable

$$\begin{aligned} g_i &= f_i, 1 \leq i \leq n-2, \\ g_{n-1} &= \frac{f_n + f_{n-1}}{2}, \\ g_n &= \frac{f_n - f_{n-1}}{2}, \end{aligned}$$

we get

$$\Phi(g_1, \dots, g_n) = \sum_{i,j=1}^n r_{i,j} g_i g_j$$

with  $r_{n,n} = 2m_{n,n-1} \neq 0$ , so we are in the situation where some diagonal entry is not zero, and we can apply the preceding transformation.

So we have decomposed

$$\Phi(f_1, \dots, f_n) = \sum_{i=1}^r c_i L_i(f_1, \dots, f_n)^2,$$

where  $r$  is the rank of  $M$ , and the  $L_i(f_1, \dots, f_n)$ 's are linearly independent linear forms, since the rank of  $M$  and the rank of the diagonal matrix with entries  $c_i$  are equal.

For the second claim, suppose that we have a second diagonal expression

$$\Phi(f_1, \dots, f_n) = \sum_{i=1}^r c'_i L'_i(f_1, \dots, f_n)^2,$$

with  $c'_i \neq 0$ , and the  $L'_i(f_1, \dots, f_n)$  are linearly independent forms, and, without loss of generality, assume that

$$\begin{aligned} c_1 > 0, \dots, c_s > 0, c_{s+1} < 0, \dots, c_r < 0, \\ c'_1 > 0, \dots, c'_{s'} > 0, c'_{s'+1} < 0, \dots, c'_r < 0, \end{aligned}$$

with  $0 \leq s \leq s' \leq r$ . If  $s' > s$ , choose values of  $f = (f_1, \dots, f_n)$  such that the values at  $f$  of the  $r - (s' - s)$  forms

$$L_1(f), \dots, L_s(f), L'_{s'+1}(f), \dots, L'_r(f)$$

are zero and the value at  $f$  of one of the forms

$$L_{s+1}(f), \dots, L_r(f)$$

is not zero.

To see that this is always possible observe that the vector subspace  $V_1$  defined by

$$L_1(f) = \dots = L_s(f) = L'_{s+1}(f) = \dots = L'_r(f) = 0$$

has dimension  $\geq n - r + s' - s > n - r$ , while the vector subspace  $V_2$  defined by

$$L_1(f) = \dots = L_s(f) = L_{s+1}(f) = \dots = L_r(f) = 0$$

has dimension  $n - r$ , since the linear forms  $L_i(f)$  are linearly independent, and thus there is a vector  $f = (f_1, \dots, f_n) \in V_1 \setminus V_2$  which satisfies

$$L_1(f) = \dots = L_s(f) = 0,$$

and  $L_i(f) \neq 0$  for some  $i$ ,  $s < i \leq r$ .

For this value of  $f = (f_1, \dots, f_n)$ ,  $\sum_{i=1}^r c_i L_i(f)^2$  is strictly negative while  $\sum_{i=1}^r c'_i L'_i(f)^2$  is non-negative. So the hypothesis  $s' > s$  leads to a contradiction. □

If  $K$  is ordered, the **signature** of  $\Phi$ ,  $\text{Sign}(\Phi)$ , is the difference between the numbers of positive  $c_i$  and negative  $c_i$  in its diagonal form.

The preceding theorem immediately implies

**Corollary 4.39.** *There exists a basis  $B$  such that, denoting also by  $B$  the matrix of  $B$  in the canonical basis,*

$$BDB^t = M$$

where  $D$  is a diagonal matrix with  $r_+$  positive entries,  $r_-$  negative entries, with  $\text{Rank}(\Phi) = r_+ + r_-$ ,  $\text{Sign}(\Phi) = r_+ - r_-$ .

Let  $R$  be a real closed field. We are going to prove that a symmetric matrix with coefficients in  $R$  can be diagonalized in a basis of orthogonal vectors.

We denote by  $u \cdot u'$  the **inner product** of vectors of  $R^n$

$$u \cdot u' = \sum_{k=1}^n u_k u'_k,$$

where  $u = (u_1, \dots, u_n)$ ,  $u' = (u'_1, \dots, u'_n)$ . The **norm** of  $u$  is  $\|u\| = \sqrt{u \cdot u}$ . Two vectors  $u$  and  $u'$  are **orthogonal** if  $u \cdot u' = 0$ .

A basis  $v_1, \dots, v_n$  of vectors of  $R^n$  is orthogonal if

$$v_i \cdot v_j = \sum_{k=1}^n v_{i,k} v_{j,k} = 0$$

for all  $i = 1, \dots, n$ ,  $j = 1, \dots, n$ ,  $j \neq i$ .

A basis  $v_1, \dots, v_n$  of vectors of  $R^n$  is **orthonormal** if it is orthogonal and moreover  $\|u\| = 1$ , for all  $i = 1, \dots, n$ .

Two linear forms

$$L = \sum_{i=1}^n u_i f_i, \quad L' = \sum_{i=1}^n u'_i f_i$$

are orthogonal if  $u \cdot u' = 0$ .

We first describe the Gram-Schmidt orthogonalization process.

**Proposition 4.40. [Gram-Schmidt orthogonalization]** *Let  $v_1, \dots, v_n$  be linearly independent vectors with coefficients in  $R$ . There is a family of linearly independent orthogonal vectors  $w_1, \dots, w_n$  with coefficients in  $R$  such that for every  $i = 1, \dots, n$ ,  $w_i - v_i$  belong to the vector space spanned by  $v_1, \dots, v_{i-1}$ .*

**Proof:** The construction proceeds by induction, starting with  $w_1 = v_1$  and continuing with

$$w_i = v_i - \sum_{j=1}^{i-1} \mu_{i,j} w_j,$$

where

$$\mu_{i,j} = \frac{v_i \cdot w_j}{\|w_j\|^2}. \quad \square$$

Let  $M$  be a symmetric matrix of dimension  $n$  with entries in  $R$ .

If  $f = (f_1, \dots, f_n)$ ,  $g = (g_1, \dots, g_n)$ , let

$$\begin{aligned} \Phi_M(f) &= f \cdot M \cdot f^t, \\ B_M(f, g) &= g \cdot M \cdot f^t, \\ u_M(f) &= M \cdot f. \end{aligned}$$



The quadratic form  $\Phi$  is **non-negative** if for every  $f \in \mathbb{R}^n$ ,  $\Phi_M(f) \geq 0$ .

**Proposition 4.41. [Cauchy-Schwarz inequality]** *If  $\Phi$  is non-negative,*

$$B_M(f, g)^2 \leq \Phi_M(f) \Phi_M(g).$$

**Proof:** Fix  $f$  and  $g$  and consider the second degree polynomial

$$P(T) = \Phi_M(f + Tg) = \Phi_M(f) + 2TB_M(f, g) + T^2\Phi_M(g).$$

For every  $t \in \mathbb{R}$ ,  $P(t)$  is non-negative since  $\Phi_M$  is non-negative. So  $P$  can be

- of degree 0 if  $\Phi_M(g) = B_M(f, g) = 0$ , in this case the inequality claimed holds
- of degree 2 with negative discriminant if  $\Phi_M(g) \neq 0$ . Since the discriminant of  $P$  is

$$4B_M(f, g)^2 - 4\Phi_M(f) \Phi_M(g),$$

the inequality claimed holds in this case too. □

Our main objective in the end of the section is to prove the following result.

**Theorem 4.42.** *Let  $M$  be a symmetric matrix with entries in  $\mathbb{R}$ . The eigenvalues of  $M$  are in  $\mathbb{R}$ , and there is an orthonormal basis of eigenvectors for  $M$  with coordinates in  $\mathbb{R}$ .*

As a consequence, since positive elements in  $\mathbb{R}$  are squares, there exists an orthogonal basis  $B$  such that, denoting also by  $B$  the matrix of  $B$  in the canonical basis,

$$B D B^t = M$$

where  $D$  is the diagonal matrix with  $r_+$  entries 1,  $r_-$  entries  $-1$ , and  $n - r$  entries 0,  $r = r_+ + r_-$ :

**Corollary 4.43.** *A quadratic form  $\Phi$  with coefficients in  $\mathbb{R}$  can always be written as*

$$\Phi = \sum_{i=1}^{r_+} L_i^2 - \sum_{i=r_++1}^{r_++r_-} L_i^2$$

where the  $L_i$  are independent orthogonal linear forms with coefficients in  $\mathbb{R}$ , and  $r = r_+ + r_-$  is the rank of  $\Phi$ .

**Corollary 4.44.** *Let  $r_+$ ,  $r_-$ , and  $r_0$  be the number of  $> 0$ ,  $< 0$ , and  $= 0$  eigenvalues of the symmetric matrix associate to the quadratic form  $\Phi$ , counted with multiplicities. Then*

$$\begin{aligned} \text{Rank}(\Phi) &= r_+ + r_-, \\ \text{Sign}(\Phi) &= r_+ - r_-. \end{aligned}$$

**Proof of Theorem 4.42:** The proof is by induction on  $n$ . The Theorem is obviously true for  $n = 1$ .

Let  $M = [m_{i,j}]_{i,j=1\dots n}$ ,  $N = [m_{i,j}]_{i,j=1\dots n-1}$ . By induction hypothesis, there exists an orthonormal matrix  $\mathcal{B}$  with entries in  $\mathbb{R}$  such that

$$\mathcal{B}^t N \mathcal{B} = D(y_1, \dots, y_{n-1})$$

where  $D(y_1, \dots, y_{n-1})$  is a diagonal matrix with entries

$$y_1 \leq \dots \leq y_{n-1}.$$

Note that the column vectors of  $\mathcal{B}$ ,  $w_1, \dots, w_{n-1}$ , form a basis of eigenvectors of the quadratic form associated to  $N$ . We can suppose without loss of generality that  $N w_i = y_i w_i$ . Let  $v_i$  be the vector of  $\mathbb{R}^n$  whose first coordinates coincide with  $w_i$  and whose last coordinate is 0 and let  $\mathcal{C}$  be an orthonormal basis completing  $v_1, \dots, v_{n-1}$  by Proposition 4.40. We have

$$\mathcal{C}^t M \mathcal{C} = \begin{bmatrix} y_1 & 0 & 0 & 0 & b_1 \\ 0 & \ddots & 0 & 0 & \vdots \\ 0 & 0 & \ddots & 0 & \vdots \\ 0 & 0 & 0 & y_{n-1} & b_{n-1} \\ b_1 & \dots & \dots & b_{n-1} & a \end{bmatrix}$$

Let  $\varepsilon$  be a variable. Define  $b'_i = b_i$  if  $b_i \neq 0$ , and  $b'_i = \varepsilon$  otherwise, and if

$$y_{i-1} < y_i = \dots = y_j < y_{j+1},$$

$y'_k = y_i + (k - i)\varepsilon$ , for  $0 \leq k \leq j - i$ . We define the symmetric matrix  $M'$  with entries in  $\mathbb{R}\langle\varepsilon\rangle$  by

$$\mathcal{C}^t M' \mathcal{C} = \begin{bmatrix} y'_1 & 0 & 0 & 0 & b'_1 \\ 0 & \ddots & 0 & 0 & \vdots \\ 0 & 0 & \ddots & 0 & \vdots \\ 0 & 0 & 0 & y'_{n-1} & b'_{n-1} \\ b'_1 & \dots & \dots & b'_{n-1} & a \end{bmatrix}.$$

Note that  $\lim_{\varepsilon} (y'_i) = y_i$ ,  $\lim_{\varepsilon} (b'_i) = b_i$ , hence  $\lim_{\varepsilon} (M') = M$ . Developing the characteristic polynomial  $P$  of  $\mathcal{C}^t M' \mathcal{C}$ , which is equal to the characteristic polynomial of  $M$ , on the last column and the last row we get

$$P = \prod_{i=1}^{n-1} (X - y'_i) (X - a') - \sum_{i=1}^{n-1} b_i^2 \prod_{j \neq i} (X - y'_j).$$

Evaluating at  $y'_i$ , we get

$$\text{sign}(P(y'_i)) = \text{sign}\left(b_i^2 \prod_{j \neq i} (y'_i - y'_j)\right) = \text{sign}(-1)^{n-i}.$$

Since the sign of  $P$  at  $-\infty$  is  $(-1)^n$ , and the sign of  $P$  at  $+\infty$  is 1, the polynomial  $P$  has  $n$  real roots satisfying

$$x'_1 < y'_1 < x'_2 < \dots < x'_{n-1} < y'_{n-1} < x'_n.$$

Taking eigenvectors of norm 1 defines an orthonormal matrix  $\mathcal{D}'$  such that

$$\mathcal{D}'^t M' \mathcal{D}' = D(x'_1, \dots, x'_n).$$

Applying  $\lim_\varepsilon$  on both sides we obtain an orthonormal matrix such that

$$\mathcal{D}^t M \mathcal{D} = D(x_1, \dots, x_n),$$

noting that  $x_1$  and  $x_n$  are bounded by an element of  $\mathbb{R}$  by Proposition 2.4. Note that  $x_1 \leq \dots \leq x_n$  are the eigenvalues of  $M$ . □

We now prove that the subdiscriminants of characteristic polynomials of symmetric matrices are sums of squares. Let  $M$  is a symmetric  $p \times p$  matrix with coefficients in a field  $K$  and  $\text{Tr}(M)$  its trace. The  $k$ -th subdiscriminant of the characteristic polynomial of  $M$   $\text{sDisc}_k(M)$  is the determinant of the matrix  $\text{Newt}_k(M)$  whose  $(i, j)$ -th entry is  $\text{Tr}(M^{i+j-2})$ ,  $i, j = 1, \dots, pk$ . Indeed, the Newton sum  $N_i$  of  $\text{CharPol}(M)$  is  $\text{Tr}(M^i)$ , the trace of the matrix  $M^i$ . If  $M$  is a symmetric  $p \times p$  matrix with coefficients in a ring  $D$ , we also define  $\text{sDisc}_k(M)$  as the determinant of the matrix  $\text{Newt}_k(M)$  whose  $(i, j)$ -th entry is  $\text{Tr}(M^{i+j-2})$ ,  $i, j = 1, \dots, p - k$ .

We define a linear basis  $E_{j,\ell}$  of the space  $\text{Sym}(p)$  of symmetric matrices of size  $p$  as follows. First define  $F_{j,\ell}$  as the matrix having all zero entries except 1 at  $(j, \ell)$ . Then take  $E_{j,j} = F_{j,j}$ ,  $E_{j,\ell} = 1/\sqrt{2}(F_{j,\ell} + F_{\ell,j})$ ,  $\ell > j$ . Define  $E$  as the ordered set  $E_{j,\ell}$   $p \geq \ell \geq j \geq 0$ , indices being taken in the order

$$(1, 1), \dots, (p, p), (1, 2), \dots, (1, p), \dots, (p - 1, p).$$

For simplicity, we index elements of  $E$  pairs  $(j, \ell)$ ,  $\ell \geq j$ .

**Proposition 4.45.** *The map associating to  $(A, B) \in \text{Sym}(p) \times \text{Sym}(p)$  the value  $\text{Tr}(AB)$  is a scalar product on  $\text{Sym}(p)$  with orthogonal basis  $E$ .*

**Proof:** Simply check. □

Let  $A_k$  be the  $(p - k) \times p(p + 1)/2$  matrix with  $(i, (j, \ell))$ -th entry the  $(j, \ell)$ -th component of  $M^{i-1}$  in the basis  $E$ .

**Proposition 4.46.**  $\text{Newt}_k(M) = A_k A_k^t$ .

**Proof:** Immediate since  $\text{Tr}(M^{i+j})$  is the scalar product of  $M^i$  by  $M^j$  in the basis  $E$ . □

We consider a generic symmetric matrix  $M = [m_{i,j}]$  whose entries are  $p(p + 1)/2$  independent variables  $m_{j,\ell}$ ,  $\ell \geq j$ . We are going to give an explicit expression of  $\text{sDisc}_k(M)$  as a sum of products of powers of 2 by squares of elements of the ring  $\mathbb{Z}[m_{j,\ell}]$ .

Let  $A_k$  be the  $(p - k) \times p(p + 1)/2$  matrix with  $(i, (j, \ell))$ -th entry the  $(j, \ell)$ -th component of  $M^{i-1}$  in the basis  $E$ .

**Theorem 4.47.**  $\text{sDisc}_k(M)$  is the sum of squares of the  $(p - k) \times (p - k)$  minors of  $A_k$ .

**Proof:** Use Proposition 4.46 and Proposition 4.10 (Cauchy-Binet formula).  $\square$

Noting that the square of a  $(p - k) \times (p - k)$  minor of  $A_k$  is a power of 2 multiplied by a square of an element of  $\mathbb{Z}[m_{j,\ell}]$ , we obtain an explicit expression of  $\text{sDisc}_k(M)$  as a sum of products of powers of 2 by squares of elements of the ring  $\mathbb{Z}[m_{j,\ell}]$ .

As a consequence the  $k$ -th subdiscriminant of the characteristic polynomial of a symmetric matrix with coefficients in a ring  $D$  is a sum of products of powers of 2 by squares of elements in  $D$ .

Let us take a simple example and consider

$$M = \begin{bmatrix} m_{11} & m_{12} \\ m_{12} & m_{22} \end{bmatrix}.$$

The characteristic polynomial of  $M$  is  $X^2 - (m_{11} + m_{22})X + m_{11}m_{22} - m_{12}^2$ , and its discriminant is  $(m_{11} + m_{22})^2 - 4(m_{11}m_{22} - m_{12}^2)$ . On the other hand the sum of the squares of the 2 by 2 minors of

$$A_0 = \begin{bmatrix} 1 & 1 & 0 \\ m_{11} & m_{22} & \sqrt{2}m_{12} \end{bmatrix}$$

is

$$(m_{22} - m_{11})^2 + (\sqrt{2}m_{12})^2 + (\sqrt{2}m_{12})^2.$$

It is easy to check the statement of Proposition 4.46 in this particular case.

**Proposition 4.48.** *Given a symmetric matrix  $M$ , there exists  $k, n - 1 \geq k \geq 0$  such that the signs of the subdiscriminants of the characteristic polynomial of  $M$  are given by*

$$\bigwedge_{p-1 \geq i \geq k} \text{sDisc}_i(M) > 0 \wedge \bigwedge_{0 \leq i < k} \text{sDisc}_i(M) = 0.$$

**Proof:** First note that, by Proposition 4.46,  $\text{sDisc}_i(M) \geq 0$ . Moreover, it follows from Proposition 4.46 that  $\text{sDisc}_i(M) = 0$  if and only if the rank of  $A_i$  is less than  $n - i$ . So,  $\text{sDisc}_{k-1}(M) = 0$  implies  $\text{sDisc}_i(M) = 0$  for every  $0 \leq i < k$  and  $\text{sDisc}_k(M) > 0$  implies  $\text{sDisc}_i(M) > 0$  for every  $n - 1 \geq i \geq k$ . In other words, for every symmetric matrix  $M$ , there exists  $k, n - 1 \geq k \geq 0$  such that the signs of the subdiscriminants of  $M$  are given by

$$\bigwedge_{p-1 \geq i \geq k} \text{sDisc}_i(M) > 0 \wedge \bigwedge_{0 \leq i < k} \text{sDisc}_i(M) = 0. \quad \square$$

As a corollary, we obtain an algebraic proof of a part of Theorem 4.42.

**Proposition 4.49.** *Let  $M$  be a symmetric matrix with entries in  $\mathbb{R}$ . The eigenvalues of  $M$  are in  $\mathbb{R}$ .*

**Proof:** The number of roots in  $\mathbb{R}$  of the characteristic polynomial  $\text{CharPol}(M)$  is  $p - k$ , using Proposition 4.48, and Theorem 4.33, while the number of distinct roots of  $\text{CharPol}(M)$  in  $\mathbb{C}$  is  $p - k$  using Proposition 4.25.  $\square$

**Proposition 4.50.** *Let  $P$  be a polynomial in  $\mathbb{R}[X]$ ,  $P = a_p X^p + \dots + a_0$ . All the roots of  $P$  are in  $\mathbb{R}$  if and only if there exists  $p > k \geq 0$  such that  $\text{sDisc}_i(P) > 0$  for all  $i$  from  $p$  to  $k$  and  $\text{sDisc}_i(P) = 0$  for all  $i$  from  $k - 1$  to  $0$*

**Proof:** Since it is clear that every polynomial having all its roots in  $\mathbb{R}$  is the characteristic polynomial of a diagonal symmetric matrix with entries in  $\mathbb{R}$ , Proposition 4.49 implies that the set of polynomials having all their roots in  $\mathbb{R}$  is contained in the set described

$$\bigvee_{k=p-1, \dots, 0} \left( \bigwedge_{p-1 \geq i \geq k} \text{sDisc}_i(P) > 0 \wedge \bigwedge_{0 \leq i < k} \text{sDisc}_i(P) = 0 \right).$$

The other inclusion follows immediately from Theorem 4.31.  $\square$

*Remark 4.51.* Note that the sign condition

$$\text{sDisc}_{p-2}(P) \geq 0 \wedge \dots \wedge \text{sDisc}_0(P) \geq 0$$

does not imply that  $P$  has all its roots in  $\mathbb{R}$ : the polynomials  $X^4 + 1$  has no real root (its four roots are  $\pm \sqrt{2}/2 \pm i\sqrt{2}/2$ , and it is immediate to check that it satisfies  $\text{sDisc}_2(P) = \text{sDisc}_1(P) = 0, \text{sDisc}_0(P) > 0$ ).

In fact, the set of polynomials having all their roots in  $\mathbb{R}$  is the closure of the set defined by

$$\text{sDisc}_{p-2}(P) > 0 \wedge \dots \wedge \text{sDisc}_0(P) > 0,$$

but does not coincide with the set defined by

$$\text{sDisc}_{p-2}(A) \geq 0 \wedge \dots \wedge \text{sDisc}_0(A) \geq 0. \quad \square$$

This is a new occurrence of the fact that the closure of a semi-algebraic set is not necessarily obtained by relaxing sign conditions defining it (see Remark 3.2).

### 4.3.2 Hermite's Quadratic Form

We define Hermite's quadratic form and indicate how its signature is related to real root counting.

Let  $\mathbb{R}$  be a real closed field,  $\mathbb{D}$  an ordered integral domain contained in  $\mathbb{R}$ ,  $\mathbb{K}$  the field of fractions of  $\mathbb{D}$ , and  $\mathbb{C} = \mathbb{R}[i]$  (with  $i^2 = -1$ ).

We consider  $P$  and  $Q$ , two polynomials in  $D[X]$ , with  $P$  monic of degree  $p$  and  $Q$  of degree  $q < p$ :

$$\begin{aligned} P &= X^p + a_{p-1}X^{p-1} + \cdots + a_1X + a_0 \\ Q &= b_qX^q + b_{q-1}X^{q-1} + \cdots + b_1X + b_0. \end{aligned}$$

We define the **Hermite quadratic form**  $\text{Her}(P, Q)$  depending of the  $p$  variables  $f_1, \dots, f_p$  in the following way:

$$\text{Her}(P, Q)(f_1, \dots, f_p) = \sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x) Q(x) (f_1 + f_2 x + \cdots + f_p x^{p-1})^2,$$

where  $\mu(x)$  is the multiplicity of  $x$ . Note that

$$\text{Her}(P, Q) = \sum_{k=1}^p \sum_{j=1}^p \sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x) Q(x) x^{k+j-2} f_k f_j.$$

When  $Q = 1$ , we get:

$$\begin{aligned} \text{Her}(P, 1) &= \sum_{k=1}^p \sum_{j=1}^p \sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x) x^{k+j-2} f_k f_j \\ &= \sum_{k=1}^p \sum_{j=1}^p N_{k+j-2} f_k f_j \end{aligned}$$

where  $N_n$  is the  $n$ -th Newton sum of  $P$  (see Definition 4.7). So the matrix associated to  $\text{Her}(P, Q)$  is  $\text{Newt}_0(P)$ .

Since the expression of  $\text{Her}(P, Q)$  is symmetric in the  $x$ 's, the quadratic form  $\text{Her}(P, Q)$  has coefficients in  $\mathbb{K}$  by Proposition 2.13. In fact, the coefficients of  $\text{Her}(P, Q)$  can be expressed in terms of the trace map.

We define  $A = \mathbb{K}[X]/(P)$ . The ring  $A$  is a  $\mathbb{K}$ -vector space of dimension  $p$  with basis  $1, X, \dots, X^{p-1}$ . Indeed any  $f \in \mathbb{K}[X]$  has a representative  $f_1 + f_2 X + \cdots + f_p X^{p-1}$  obtained by taking its remainder in the euclidean division by  $P$ , and if  $f$  and  $g$  are equal modulo  $P$ , their remainder in the euclidean division by  $P$  are equal.

We denote by  $\text{Tr}$  the usual **trace** of a linear map from a finite dimensional vector space  $A$  to  $A$ , which is the sum of the entries on the diagonal of its associated matrix in any basis of  $A$ .

**Notation 4.52. [Multiplication map]** For  $f \in A$ , we denote by  $L_f: A \rightarrow A$  the linear map of multiplication by  $f$ , sending any  $g \in A$  to the remainder of  $fg$  in the euclidean division by  $P$ .  $\square$

**Proposition 4.53.** *The quadratic form  $\text{Her}(P, Q)$  is the quadratic form associating to*

$$f = f_1 + f_2 X + \cdots + f_p X^{p-1} \in A = \mathbb{K}[X]/(P)$$

*the expression  $\text{Tr}(L_Q f^2)$ .*

The proof of Proposition 4.53 relies on the following results.

**Proposition 4.54.**

$$\mathrm{Tr}(L_f) = \sum_{x \in \mathrm{Zer}(P, C)} \mu(x) f(x).$$

**Proof:** The proof proceeds by induction on the number of distinct roots of  $P$ .

When  $P = (X - x)^{\mu(x)}$ , since  $x$  is root of  $f - f(x)$ ,

$$(f - f(x))^{\mu(x)} = 0 \text{ modulo } P$$

and  $L_{f-f(x)}$  is nilpotent, with characteristic polynomial  $X^{\mu(x)}$ . Thus  $L_{f-f(x)}$  has a unique eigenvalue 0 with multiplicity  $\mu(x)$ . So  $\mathrm{Tr}(L_{f-f(x)}) = 0$  and  $\mathrm{Tr}(L_f) = \mu(x) f(x)$ .

If  $P = P_1 P_2$  with  $P_1$  and  $P_2$  coprime, by Proposition 1.9 there exists  $U_1$  and  $U_2$  with  $U_1 P_1 + U_2 P_2 = 1$ . Let

$$e_1 = U_2 P_2 = 1 - U_1 P_1, \quad e_2 = U_1 P_1 = 1 - U_2 P_2.$$

It is easy to verify that

$$e_1^2 = e_1, \quad e_2^2 = e_2, \quad e_1 e_2 = 0, \quad e_1 + e_2 = 1$$

in A. It is also easy to check that the mapping from  $\mathbb{K}[X]/(P_1) \times \mathbb{K}[X]/(P_2)$  to  $\mathbb{K}[X]/(P)$  associating to  $(Q_1, Q_2)$  the polynomial  $Q = Q_1 e_1 + Q_2 e_2$  is an isomorphism. Moreover, if  $f_1 = f \bmod P_1$  and  $f_2 = f \bmod P_2$ ,  $\mathbb{K}[X]/(P_1)$  and  $\mathbb{K}[X]/(P_2)$  are stable by  $L_f$  and  $L_{f_1}$  and  $L_{f_2}$  are the restrictions of  $L_f$  to  $\mathbb{K}[X]/(P_1)$  and  $\mathbb{K}[X]/(P_2)$ . Then  $\mathrm{Tr}(L_f) = \mathrm{Tr}(L_{f_1}) + \mathrm{Tr}(L_{f_2})$ . This proves the proposition by induction, since the number of roots of  $P_1$  and  $P_2$  are smaller than the number of roots of  $P$ .  $\square$

**Proposition 4.55.** *Let  $C = \mathrm{Quo}(P'Q, P)$ , then*

$$\frac{P'Q}{P} = C + \sum_{n=0}^{\infty} \frac{\mathrm{Tr}(L_Q X^n)}{X^{n+1}}.$$

**Proof:** As already seen in the proof of Proposition 4.8

$$\frac{P'}{P} = \sum_{x \in \mathrm{Zer}(P, C)} \frac{\mu(x)}{(X-x)}.$$

Dividing  $Q$  by  $X - x$  and letting  $C_x$  be the quotient,

$$Q = Q(x) + (X - x)C_x,$$

and thus

$$\frac{P'Q}{P} = \sum_{x \in \mathrm{Zer}(P, C)} \mu(x) \left( C_x + \frac{Q(x)}{(X-x)} \right).$$

Since

$$\frac{1}{X-x} = \sum_{n=0}^{\infty} \frac{x^n}{X^{n+1}},$$

the coefficient of  $1/X^{n+1}$  in the development of  $P'Q/P$  in powers of  $1/X$  is thus,

$$\sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x) Q(x) x^n.$$

Now apply Proposition 4.54 □

**Proof of Proposition 4.53:** By Proposition 4.55,

$$\text{Tr}(L_{QX^{k+j}}) = \sum_{x \in \text{Zer}(P, \mathbb{C})} \mu(x) Q(x) x^{k+j}.$$

In other words,  $\text{Tr}(L_{QX^{k+j}})$  is the  $j + 1, k + 1$ -th entry of the symmetric matrix associated to Hermite's quadratic form  $\text{Her}(P, Q)$  in the basis  $1, X, \dots, X^{p-1}$ . □

Note that Proposition 4.55 implies that the coefficients of  $\text{Her}(P, Q)$  belong to  $\mathbb{D}$ , since  $L_f$  expressed in the canonical basis has entries in  $\mathbb{D}$ .

*Remark 4.56.* As a consequence of Proposition 4.53, the quadratic form  $\text{Her}(P, 1)$  is the quadratic form associating to

$$f = f_1 + f_2 X \dots + f_p X^{p-1} \in \mathbb{A} = \mathbb{K}[X]/(P)$$

the expression  $\text{Tr}(L_{f^2})$ . So the  $j + 1, k + 1$ -th entry of the symmetric matrix associated to Hermite's quadratic form  $\text{Her}(P, 1)$  in the basis  $1, X, \dots, X^{p-1}$  is  $\text{Tr}(L_{X^{j+k}}) = N_{k+j}$ . Note that Proposition 4.55 is a generalization of Proposition 4.8. □

The main result about Hermite's quadratic form is the following theorem. We use again the notation

$$\text{TaQ}(Q, P) = \sum_{x \in \mathbb{R}, P(x)=0} \text{sign}(Q(x)).$$

**Theorem 4.57. [Hermite]**

$$\begin{aligned} \text{Rank}(\text{Her}(P, Q)) &= \#\{x \in \mathbb{C} \mid P(x) = 0 \wedge Q(x) \neq 0\}, \\ \text{Sign}(\text{Her}(P, Q)) &= \text{TaQ}(Q, P). \end{aligned}$$

As an immediate consequence

**Theorem 4.58.** *The rank of  $\text{Her}(P, 1)$  is equal to the number of roots of  $P$  in  $\mathbb{C}$ . The signature of  $\text{Her}(P, 1)$  is equal to the number of roots of  $P$  in  $\mathbb{R}$ .*



**Proof of Theorem 4.57:** For  $x \in \mathbb{C}$ , let  $L(x, -)$  be the linear form on  $\mathbb{C}^p$  defined by:

$$L(x, f) = f_1 + f_2 x + \cdots + f_p x^{p-1}.$$

Let  $\{x \in \mathbb{C} \mid P(x) = 0 \wedge Q(x) \neq 0\} = \{x_1, \dots, x_r\}$ . Thus,

$$\text{Her}(P, Q) = \sum_{i=1}^r \mu(x_i) Q(x_i) L(x_i, f)^2.$$

The linear forms  $L(x_i, f)$  are linearly independent since the roots are distinct and the Vandermonde determinant

$$\det(V(x_1, \dots, x_r)) = \prod_{r \geq i > j \geq 1} (x_i - x_j).$$

is non-zero. Thus the rank of  $\text{Her}(P, Q)$  is  $r$ .

Let

$$\{x \in \mathbb{R} \mid P(x) = 0 \wedge Q(x) \neq 0\} = \{y_1, \dots, y_s\}.$$

$$\{x \in \mathbb{C} \setminus \mathbb{R} \mid P(x) = 0 \wedge Q(x) \neq 0\} = \{z_1, \bar{z}_1, \dots, z_t, \bar{z}_t\}.$$

The quadratic form  $\text{Her}(P, Q)$  is equal to

$$\sum_{i=1}^s \mu(y_i) Q(y_i) L(y_i, f)^2 + \sum_{j=1}^t \mu(z_j) (Q(z_j) L(z_j, f)^2 + Q(\bar{z}_j) L(\bar{z}_j, f)^2),$$

with the  $L(y_i, f)$ ,  $L(z_j, f)$ ,  $L(\bar{z}_j, f)$  ( $i = 1, \dots, s$ ,  $j = 1, \dots, t$ ) linearly independent.

Writing  $\mu(z_j) Q(z_j) = (a(z_j) + i b(z_j))^2$  with  $a(z_j), b(z_j) \in \mathbb{R}$  and denoting by  $s_i(z_j)$  and  $t_i(z_j)$  the real and imaginary part of  $z_j^i$ ,

$$L_1(z_j) = \sum_{i=1}^p (a(z_j) s_i(z_j) - b(z_j) t_i(z_j)) f_i$$

$$L_2(z_j) = \sum_{i=1}^p (a(z_j) t_i(z_j) + b(z_j) s_i(z_j)) f_i$$

are linear forms with coefficients in  $\mathbb{R}$  such that

$$\mu(z_j) (Q(z_j) L(z_j, f)^2 + Q(\bar{z}_j) L(\bar{z}_j, f)^2) = 2L_1(z_j)^2 - 2L_2(z_j)^2.$$

Moreover the  $L(y_i, f)$ ,  $L_1(z_j)$ ,  $L_2(z_j)$  ( $i = 1, \dots, s$ ,  $j = 1, \dots, t$ ) are linearly independent linear forms. So, using Theorem 4.38 (Sylvester's inertia law), the signature of  $\text{Her}(P, Q)$  is the signature of  $\sum_{i=1}^s \mu(y_i) Q(y_i) L(y_i, f)^2$ . Since the linear forms  $L(y_i, f)$  are linearly independent, the signature of  $\text{Her}(P, Q)$  is  $\text{TaQ}(Q, P)$ .  $\square$

*Remark 4.59.* Note that it follows from Theorem 4.58 and Theorem 4.33 that the signature of  $\text{Her}(P, 1)$ , which is the number of roots of  $P$  in  $\mathbb{R}$  can be computed from the signs of the principal minors  $s\text{Disc}_p(P), k = 1, \dots, p$  of the symmetric matrix  $\text{Newt}_0(P)$  defining  $\text{Her}(P, 1)$ . This is a general fact about Hankel matrices that we shall define and study in Chapter 9.  $\square$

## 4.4 Polynomial Ideals

### 4.4.1 Hilbert's Basis Theorem

An **ideal**  $I$  of a ring  $A$  is a subset  $I \subset A$  containing 0 that is closed under addition and under multiplication by any element of  $A$ . To an ideal  $I$  of  $A$  is associated an equivalence relation on  $A$  called **congruence modulo**  $I$ . We write  $a = b \pmod I$  if and only if  $a - b \in I$ . It is clear that if  $a_1 - b_1 \in I, a_2 - b_2 \in I$  then  $(a_1 + a_2) - (b_1 + b_2) \in I, a_1 a_2 - b_1 b_2 = a_1(a_2 - b_2) + b_2(a_1 - b_1) \in I$ .

The **quotient ring**  $A/I$  is the set of equivalence classes equipped with the natural ring structure obtained by defining the sum or product of two classes as the class of the sum or product of any members of the classes. Observation 4.9 shows that this is well defined.

The set of those elements  $a$  such that a power of  $a$  belongs to the ideal  $I$  is an ideal called the **radical of**  $I$ :

$$\sqrt{I} = \{a \in A \mid \exists m \in \mathbb{N} \quad a^m \in I\}.$$

A **prime ideal** is an ideal such that  $xy \in I$  implies  $x \in I$  or  $y \in I$ .

To a finite set of polynomials  $\mathcal{P} \subset \mathbb{K}[X_1, \dots, X_k]$  is associated  $\text{Ideal}(\mathcal{P}, \mathbb{K})$ , the **ideal generated by**  $\mathcal{P}$  in  $\mathbb{K}[X_1, \dots, X_k]$  i.e.,

$$\text{Ideal}(\mathcal{P}, \mathbb{K}) = \left\{ \sum_{P \in \mathcal{P}} A_P P \mid A_P \in \mathbb{K}[X_1, \dots, X_k] \right\}.$$

A polynomial in  $\text{Ideal}(\mathcal{P}, \mathbb{K})$  vanishes at every point of  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .

Note that when  $k = 1$ , the ideal generated by  $\mathcal{P}$  in  $\mathbb{K}[X_1]$  is **principal** (i.e. generated by a single polynomial) and generated by the greatest common divisor of the polynomials in  $\mathcal{P}$  (Definition, page 13).

This is no longer true for a general  $k$ , but the following finiteness theorem holds.

**Theorem 4.60. [Hilbert's basis theorem]** *Any ideal  $I \subset \mathbb{K}[X_1, \dots, X_k]$  is finitely generated, i.e. there exists a finite set  $\mathcal{P}$  such that  $I = \text{Ideal}(\mathcal{P}, \mathbb{K})$ .*

The proof uses the partial order of divisibility on the set  $\mathcal{M}_k$  of monomials in  $k$  variables  $X_1, \dots, X_k$ , which can be identified with  $\mathbb{N}^k$ , partially ordered by

$$\alpha = (\alpha_1, \dots, \alpha_k) \prec \beta = (\beta_1, \dots, \beta_k) \Leftrightarrow \alpha_1 \leq \beta_1, \dots, \alpha_k \leq \beta_k.$$

If  $\alpha = (\alpha_1, \dots, \alpha_{k-1}) \in \mathbb{N}^{k-1}$  and  $n \in \mathbb{N}$ , we denote by  $(\alpha, n) = (\alpha_1, \dots, \alpha_{k-1}, n)$ .

**Lemma 4.61. [Dickson’s lemma]** *Every subset of  $\mathcal{M}_k$  closed under multiplication has a finite number of minimal elements with respect to the partial order of divisibility.*

**Proof:** The proof is by induction on the number  $k$  of variables. If  $k = 1$ , the result is clear. Suppose that the property holds for  $k - 1$ . Let  $B \subset \mathcal{M}_k$  and

$$A = \left\{ X^\alpha \in \mathcal{M}_{k-1} \mid \exists n \in \mathbb{N} X^{(\alpha, n)} \in B \right\}.$$

By induction hypothesis,  $A$  has a finite set of minimal elements for the partial order of divisibility

$$\left\{ X^{\alpha(1)}, \dots, X^{\alpha(N)} \right\}.$$

Let  $n$  be such that for every  $i = 1, \dots, N$ ,  $X^{(\alpha(i), n)} \in B$ . For every  $m < n$ ,

$$C_m = \left\{ X^\alpha \in \mathcal{M}_{k-1} \mid X^{(\alpha, m)} \in B \right\}$$

has a finite set of minimal elements with respect to the partial order of divisibility

$$\left\{ X^{\gamma(m, 1)}, \dots, X^{\gamma(m, \ell(m))} \right\},$$

using again the induction hypothesis. Consider the finite set

$$D = \left\{ X^{(\alpha(i), n)} \mid i = 1, \dots, N \right\} \bigcup_{m=0}^n \left\{ X^{(\gamma(m, i), m)} \mid i = 1, \dots, \ell(m) \right\}.$$

Let  $X^\beta \in B$ , with  $\beta = (\alpha, r)$ . If  $r \geq n$ ,  $X^\beta$  is multiple of  $X^{(\alpha(i), n)}$  for some  $i = 1, \dots, N$ . On the other hand, if  $r < n$ ,  $X^\beta$  is multiple of  $X^{(\gamma(r, i), r)}$  for some  $i = 1, \dots, \ell(r)$ . So every element of  $B$  is multiple of an element in  $D$ . It is clear that a finite number of minimal elements for the partial order of divisibility can be extracted from  $D$ . □

In order to prove Theorem 4.60, the notion of monomial ordering is useful.

**Definition 4.62. [Monomial ordering]** A total ordering on the set  $\mathcal{M}_k$  of monomials in  $k$  variables is a **monomial ordering** if the following properties hold

- a)  $X^\alpha > 1$  for every  $\alpha \in \mathbb{N}^k$ ,  $\alpha \neq (0, \dots, 0)$
- b)  $X_1 > \dots > X_k$ ,
- c)  $X^\alpha > X^\beta \implies X^{\alpha+\gamma} > X^{\beta+\gamma}$ , for every  $\alpha, \beta, \gamma$  elements of  $\mathbb{N}^k$ ,
- d) every decreasing sequence of monomials for the monomial order  $<$  is finite. □

The lexicographical ordering defined in Notation 2.14 and the graded lexicographical ordering defined in Notation 2.15 are examples of monomial orderings. Another important example of monomial ordering is the reverse lexicographical ordering defined above.

**Definition 4.63. [Reserve lexicographical ordering]** The **reverse lexicographical ordering**,  $<_{\text{revlex}}$ , on the set  $\mathcal{M}_k$  of monomials in  $k$  variables is the total order  $X^\alpha <_{\text{grlex}} X^\beta$  defined by

$$X^\alpha <_{\text{grlex}} X^\beta \Leftrightarrow (\deg(X^\alpha) < \deg(X^\beta)) \vee (\deg(X^\alpha) = \deg(X^\beta) \wedge \bar{\beta} <_{\text{lex}} \bar{\alpha})$$

with  $\alpha = (\alpha_1, \dots, \alpha_k)$ ,  $\beta = (\beta_1, \dots, \beta_k)$ ,  $\bar{\alpha} = (\alpha_k, \dots, \alpha_1)$ ,  $\bar{\beta} = (\beta_k, \dots, \beta_1)$ ,  $X^\alpha = X_1^{\alpha_1} \dots X_k^{\alpha_k}$ ,  $X^\beta = X_1^{\beta_1} \dots X_k^{\beta_k}$ , and  $<_{\text{lex}}$  is the lexicographical ordering defined in Notation 2.14.

In the reverse lexicographical ordering above,  $X_1 >_{\text{revlex}} \dots >_{\text{revlex}} X_k$ . The smallest monomial with respect to the reverse lexicographical ordering is 1, and the reverse lexicographical ordering order is compatible with multiplication. Note that the set of monomials less than or equal to a monomial  $X^\alpha$  in the reverse lexicographical ordering is finite.  $\square$

**Definition 4.64.** Given a polynomial  $P \in K[X_1, \dots, X_k]$  we write  $\text{cof}(X^\alpha, P)$  for the coefficient of the monomial  $X^\alpha$  in the polynomial  $P$ . The monomial  $X^\alpha$  is a **monomial of**  $P$  if  $\text{cof}(X^\alpha, P) \neq 0$ , and  $\text{cof}(X^\alpha, P)X^\alpha$  is a **term of**  $P$ .

Given a monomial ordering  $<$  on  $\mathcal{M}_k$ , we write  $\text{lmon}(P)$  for the **leading monomial** of  $P$  with respect to  $<$  i.e. the largest monomial of  $P$  with respect to  $<$ . The **leading coefficient** of  $P$  is  $\text{lcof}(P) = \text{cof}(\text{lmon}(P), P)$ , and the **leading term** of  $P$  is  $\text{lt}(P) = \text{lcof}(P)\text{lmon}(P)$ . Let  $X^\alpha$  be a monomial of  $P$ , and let  $G$  be another polynomial. The **reduction** of  $(P, X^\alpha)$  by  $G$  is defined by

$$\begin{aligned} & \text{Red}(P, X^\alpha, G) \\ = & \begin{cases} P - (\text{cof}(X^\alpha, P)/\text{lcof}(G)) X^\beta G & \text{if } \exists \beta \in \mathbb{N}^k \ X^\alpha = X^\beta \text{lmon}(G), \\ P & \text{otherwise.} \end{cases} \end{aligned}$$

Given a finite set of polynomials,  $\mathcal{G} \subset K[X_1, \dots, X_k]$ ,  $Q$  is a **reduction** of  $P$  modulo  $\mathcal{G}$  if there is a  $G \in \mathcal{G}$  and a monomial  $X^\alpha$  of  $P$  such that  $Q = \text{Red}(P, X^\alpha, G)$ . We say that  $P$  is **reducible** to  $Q$  modulo  $\mathcal{G}$  if there is a finite sequence of reductions modulo  $\mathcal{G}$  starting with  $P$  and ending at  $Q$ .  $\square$

*Remark 4.65.* Note that if  $P$  is reducible to  $Q$  modulo  $\mathcal{G}$ , it follows that  $(P - Q) \in \text{Ideal}(\mathcal{G}, K)$ . Note also that if  $P$  is reducible to 0 modulo  $\mathcal{G}$ , then

$$\exists G_1 \in \mathcal{G} \dots \exists G_s \in \mathcal{G} \ P = A_1 G_1 + \dots + A_s G_s,$$

with  $\text{lmon}(A_i G_i) \leq \text{lmon}(P)$  for all  $i = 1, \dots, s$ .  $\square$

**Definition 4.66.** A **Gröbner basis** of an ideal  $I \subset K[X_1, \dots, X_k]$  for the monomial ordering  $<$  on  $\mathcal{M}_k$  is a finite set,  $\mathcal{G} \subset I$ , such that

- the leading monomial of any element in  $I$  is a multiple of the leading monomial of some element in  $\mathcal{G}$ ,
- the leading monomial of any element of  $\mathcal{G}$  is not a multiple of the leading monomial of another element in  $\mathcal{G}$ .

A **Gröbner basis** for the monomial ordering  $<$  on  $\mathcal{M}_k$  is a finite set  $\mathcal{G} \subset K[X_1, \dots, X_k]$  which is a Gröbner basis of the ideal  $\text{Ideal}(\mathcal{G}, K)$ . □

Deciding whether an element belongs to an ideal  $I$  is easy given a Gröbner basis of  $I$ .

**Proposition 4.67.** *If  $\mathcal{G}$  is a Gröbner basis of  $I$  for the monomial ordering  $<$  on  $\mathcal{M}_k$ ,  $P \in I$  if and only if  $P$  is reducible to 0 modulo  $\mathcal{G}$ .*

**Proof:** It is clear that if  $P$  is reducible to 0 modulo  $\mathcal{G}$ ,  $P \in I$ . Conversely, let  $P \neq 0 \in I$ . Then, the leading monomial of  $P$  is a multiple of the leading monomial of some  $G \in \mathcal{G}$ , so that defining  $Q = \text{Red}(P, \text{lmon}(P), G)$  either  $Q = 0$  or  $\text{lmon}(Q) < \text{lmon}(P)$ ,  $Q \in I$ . Since there is non infinite decreasing sequence for  $<$ , this process must terminate at zero after a finite number of steps. □

As a consequence

**Proposition 4.68.** *A Gröbner basis of  $I$  for the monomial ordering  $<$  on  $\mathcal{M}_k$  is a set of generators of  $I$ .*

**Proof:** Let  $P \in I$ . By Proposition 4.67,  $P$  is reducible to 0 by  $\mathcal{G}$  and  $P \in I(\mathcal{G}, K)$ . □

**Proposition 4.69.** *Every ideal of  $K[X_1, \dots, X_k]$  has a Gröbner basis for any monomial ordering  $<$  on  $\mathcal{M}_k$ .*

**Proof:** Let  $I \subset K[X_1, \dots, X_k]$  be an ideal and let  $\text{lmon}(I)$  be the set of leading monomials of elements of  $I$ . By Lemma 4.61, there is a finite set of minimal elements in  $\text{lmon}(I)$  for the partial order of divisibility, denoted by  $\{X^{\alpha(1)}, \dots, X^{\alpha(N)}\}$ . Let  $\mathcal{G} = \{G_1, \dots, G_N\}$  be elements of  $I$  with leading monomials  $\{X^{\alpha(1)}, \dots, X^{\alpha(N)}\}$ . By definition of  $\mathcal{G}$ , the leading monomial of any polynomial in  $I$  is a multiple of the leading monomial of some polynomial in  $\mathcal{G}$ , and no leading monomial of  $\mathcal{G}$  is divisible by another leading monomial of  $\mathcal{G}$ . □

**Proof of Theorem 4.60:** The claim is an immediate corollary of Proposition 4.69 since a Gröbner basis of an ideal is a finite number of generators, by Proposition 4.68. □

**Corollary 4.70.** *Let  $I_1 \subset I_2 \subset \dots \subset I_n \subset \dots$  be an ascending chain of ideals of  $K[X_1, \dots, X_k]$ . Then  $\exists n \in \mathbb{N} \forall m \in \mathbb{N} (m > n \Rightarrow I_m = I_n)$ .*

**Proof:** It is clear that  $I = \bigcup_{i \geq 0} I_i$  is an ideal and has a finite set of generators according to Theorem 4.60. This finite set of generators belongs to some  $I_N$  and so  $I_N = I$ .  $\square$

If  $I \subset K[X_1, \dots, X_k]$  is an ideal and  $L$  is a field containing  $K$ , we denote by  $\text{Zer}(I, L^k)$  the set of common zeros of  $I$  in  $L^k$ ,

$$\text{Zer}(I, L^k) = \{x \in L^k \mid \forall P \in I \ P(x) = 0\}.$$

When  $L = K$ , this defines the **algebraic sets** contained in  $K^k$ . Note that Theorem 4.60 implies that every algebraic set contained in  $K^k$  is of the form

$$\text{Zer}(\mathcal{P}, K^k) = \{x \in K^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\},$$

where  $\mathcal{P}$  is a finite set of polynomials, so that the definition of algebraic sets given here coincides with the definition of algebraic sets given in Chapter 1 (Definition page 11) when  $K = \mathbb{C}$  and in Chapter 2 (Definition page 57) when  $K = \mathbb{R}$ .

#### 4.4.2 Hilbert's Nullstellensatz

Hilbert's Nullstellensatz (weak form) is the following result.

**Theorem 4.71. [Weak Hilbert's Nullstellensatz]** *Let  $\mathcal{P} = \{P_1, \dots, P_s\}$  be a finite subset of  $K[X_1, \dots, X_k]$  then  $\text{Zer}(\mathcal{P}, K^k) = \emptyset$  if and only if there exist  $A_1, \dots, A_s \in K[X_1, \dots, X_k]$  such that*

$$A_1 P_1 + \dots + A_s P_s = 1.$$

We develop several tools and technical results before proving it.

The **degree** of a monomial  $X^\alpha = X_k^{\alpha_k} \dots X_1^{\alpha_1}$  in  $k$  variables is the sum of the degrees with respect to each variable and the **degree** of a polynomial  $P$  in  $k$  variables, denoted  $\text{deg}(Q)$ , is the maximum degree of its monomials. A polynomial is **homogeneous** if all its monomials have the same degree.

**Definition 4.72.** A non-zero polynomial  $P \in K[X_1, \dots, X_{k-1}][X_k]$  is **quasi-monic** with respect to  $X_k$  if its leading coefficient with respect to  $X_k$  is an element of  $K$ . A set of polynomials  $\mathcal{P}$  is quasi-monic with respect to  $X_k$  if each polynomial in  $\mathcal{P}$  is quasi-monic with respect to  $X_k$ .  $\square$

If  $v$  is a linear automorphism  $K^k \rightarrow K^k$  and  $[v_{i,j}]$  is its matrix in the canonical basis, we write

$$v(X) = \left( \sum_{j=1}^k v_{1,j} X_j, \dots, \sum_{j=1}^k v_{k,j} X_j \right).$$

**Lemma 4.73.** *Let  $\mathcal{P} \subset \mathbb{K}[X_1, \dots, X_k]$  be a finite subset. Then, there exists a linear automorphism  $v: \mathbb{K}^k \rightarrow \mathbb{K}^k$  such that for all  $P \in \mathcal{P}$ , the polynomial  $P(v(X))$  is quasi-monic in  $X_k$ .*

**Proof:** Choose a linear automorphism of the form

$$v(X_1, \dots, X_k) = (X_1 + a_1 X_k, X_2 + a_2 X_k, \dots, X_{k-1} + a_{k-1} X_k, X_k)$$

with  $a_i \in \mathbb{K}$ . Writing  $P(X) = \Pi(X) + \dots$ , where  $\Pi$  is the homogeneous part of highest degree (say  $d$ ) of  $P$ , we have

$$P(v(X)) = \Pi(a_1, \dots, a_{k-1}, 1) X_k^d + Q$$

(where  $Q$  has smaller degree in  $X_k$ ); it is enough to choose  $a_1, \dots, a_{k-1}$  such that none of the  $\Pi(a_1, \dots, a_{k-1}, 1)$  is zero. This can be done by taking the product of the  $\Pi$  and using the following Lemma 4.74. □

**Lemma 4.74.** *If a polynomial  $B(Z_1, \dots, Z_k)$  in  $\mathbb{K}[Z_1, \dots, Z_k]$  is not identically zero and has degree  $d$ , there are elements  $(z_1, \dots, z_k)$  in  $\{0, \dots, d\}^k$  such that  $B(z_1, \dots, z_k)$  is a non-zero element of  $\mathbb{K}$ .*

**Proof:** The proof is by induction on  $k$ . It is true for a polynomial in one variable since a non-zero polynomial of degree  $d$  has at most  $d$  roots in a field, so it does not vanish on at least one point of  $\{0, \dots, d\}$ . Suppose now that it is true for  $k - 1$  variables, and consider a polynomial  $B(Z_1, \dots, Z_k)$  in  $k$  variables of degree  $d$  that is not identically zero. Thus, if we consider  $B$  as a polynomial in  $Z_k$  with coefficients in  $\mathbb{K}[Z_1, \dots, Z_{k-1}]$ , one of its coefficients is not identically zero in  $\mathbb{K}[Z_1, \dots, Z_{k-1}]$ . Hence, by the induction hypothesis, there exist  $(z_1, \dots, z_{k-1})$  in  $\{0, \dots, d\}^{k-1}$  with  $B(z_1, \dots, z_{k-1}, Z_k)$  not identically zero. The degree of  $B(z_1, \dots, z_{k-1}, Z_k)$  is at most  $d$ , so we have reduced to the case of one variable, which we have already considered. □

Let  $\mathcal{P} \subset \mathbb{K}[X_1, \dots, X_k]$  and  $\mathcal{Q} \subset \mathbb{K}[X_1, \dots, X_{k-1}]$  be two finite sets of polynomials. The projection  $\pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k-1}$  forgetting the last coordinate is a **finite mapping from  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  onto  $\text{Zer}(\mathcal{Q}, \mathbb{C}^{k-1})$**  if its restriction to  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  is surjective on  $\text{Zer}(\mathcal{Q}, \mathbb{C}^{k-1})$  and if  $\mathcal{P}$  contains a polynomial quasi-monic in  $X_k$ , denoted by  $P$ . Since  $P$  is quasi-monic in  $X_k$ , for every  $y \in \text{Zer}(\mathcal{Q}, \mathbb{C}^{k-1})$ ,  $\text{Zer}(P(y, X_k), \mathbb{C})$  is finite. Thus

$$\pi^{-1}(y) \cap \text{Zer}(\mathcal{P}, \mathbb{C}^k) \subset \text{Zer}(P(y, X_k), \mathbb{C})$$

is finite.

**Proposition 4.75.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{K}[X_1, \dots, X_k]$  with  $P_1$  quasi-monic in  $X_k$ . There exists a finite set*

$$\text{Proj}_{X_k}(\mathcal{P}) \subset \mathbb{K}[X_1, \dots, X_{k-1}] \cap \text{Ideal}(\mathcal{P}, \mathbb{K})$$

such that

$$\pi(\text{Zer}(\mathcal{P}, \mathbb{C}^k)) = \text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1})$$

(where  $\pi$  is the projection from  $\mathbb{C}^k$  to  $\mathbb{C}^{k-1}$  forgetting the last coordinate) and  $\pi$  is a finite mapping from  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  to  $\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1})$ .

**Proof:** If  $s = 1$ , take  $\text{Proj}_{X_k}(\mathcal{P}) = \{0\}$ . Since  $P_1$  is quasi-monic, the conclusion is clear.

If  $s > 1$ , then let  $U$  be a new indeterminate, and let

$$R(U, X_1, \dots, X_k) = P_2(X) + UP_3(X) + \dots + U^{s-2}P_s(X).$$

The resultant of  $P_1$  and  $R$  with respect to  $X_k$  (apply definition in page 106), belongs to  $\mathbb{K}[U, X_1, \dots, X_{k-1}]$  and is written

$$\text{Res}_{X_k}(P_1, R) = Q_t U^{t-1} + \dots + Q_1,$$

with  $Q_i \in \mathbb{K}[X_1, \dots, X_{k-1}]$ . Let  $\text{Proj}_{X_k}(\mathcal{P}) = \{Q_1, \dots, Q_t\}$ .

By Proposition 4.18, there are polynomials  $M$  and  $N$  in  $\mathbb{K}[U, X_1, \dots, X_k]$  such that

$$\text{Res}_{X_k}(P_1, R) = MP_1 + NR.$$

Identifying the coefficients of the powers of  $U$  in this equality, one sees that there are for  $i = 1, \dots, t$  identities  $Q_i = M_i P_1 + N_{i,2} P_2 + \dots + N_{i,s} P_s$  with  $M_i$  and  $N_{i,2}, \dots, N_{i,s}$  in  $\mathbb{K}[X_1, \dots, X_k]$  so that  $Q_1, \dots, Q_t$  belong to  $\text{Ideal}(\mathcal{P}, \mathbb{K}) \cap \mathbb{K}[X_1, \dots, X_{k-1}]$ .

Since

$$\text{Proj}_{X_k}(\mathcal{P}) \subset \text{Ideal}(\mathcal{P}, \mathbb{K}) \cap \mathbb{K}[X_1, \dots, X_{k-1}],$$

it follows that

$$\pi(\text{Zer}(\mathcal{P}, \mathbb{C}^k)) \subset \text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1}).$$

In the other direction, suppose  $x' \in \text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1})$ . Then for every  $u \in \mathbb{C}$ , we have  $\text{Res}_{X_k}(P_1, R)(u, x') = 0$ . Since  $P_1$  is quasi-monic with respect to  $X_k$ ,

$$\text{Res}(P_1(x', X_k), R(u, x', X_k)) = \text{Res}_{X_k}(P_1, R)(u, x') = 0,$$

using Proposition 4.20. For every  $u \in \mathbb{C}$ , by Proposition 4.15 the polynomials  $P(x', X_k)$  and  $R(u, x', X_k)$  have a common factor in  $\mathbb{K}[X_k]$ , hence a common root in  $\mathbb{C}$ . Since  $P(x', X_k)$  has a finite number of roots in  $\mathbb{C}$ , one of them, say  $x_k$ , is a root of  $R(u, x', X_k)$  for infinitely many  $u \in \mathbb{C}$ . Choosing  $s - 1$  such distinct elements  $u_1, \dots, u_{s-1}$ , we get that the polynomial  $R(U, x', x_k)$  of degree  $\leq s - 2$  in  $U$  has  $s - 1$  distinct roots, which is possible only if  $R(U, x', x_k)$  is identically zero. So one has  $P_2(x', x_k) = \dots = P_s(x', x_k) = 0$ . We have proved that for any  $x' \in \text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1})$ , there exist a finite non-zero number of  $x_k$  such that  $(x', x_k) \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , so that

$$\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1}) \subset \pi(\text{Zer}(\mathcal{P}, \mathbb{C}^k)).$$

Since  $P_1$  is monic,  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  is finite on  $\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}), \mathbb{C}^{k-1})$ .  $\square$



Let  $\mathcal{P} \subset \mathbb{K}[X_1, \dots, X_k]$ ,  $\mathcal{T} \subset \mathbb{K}[X_1, \dots, X_{k'}]$  be two finite sets of polynomials and  $k > k'$ . The projection  $\Pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k'}$  forgetting the last  $(k - k')$  coordinates is a **finite mapping from  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  to  $\text{Zer}(\mathcal{T}, \mathbb{C}^{k'})$**  if for each  $i$ ,  $0 \leq i \leq k - k'$  there exists a finite set of polynomials  $\mathcal{Q}_{k-i} \subset \mathbb{K}[X_1, \dots, X_{k-i}]$  with  $\mathcal{P} = \mathcal{Q}_k, \mathcal{T} = \mathcal{Q}_{k'}$  such that for every  $i$ ,  $0 \leq i \leq k - k' - 1$ , the projection from  $\mathbb{C}^{k-i}$  to  $\mathbb{C}^{k-i-1}$  forgetting the last coordinate is a finite mapping from  $\text{Zer}(\mathcal{Q}_{k-i}, \mathbb{C}^{k-i})$  to  $\text{Zer}(\mathcal{Q}_{k-i-1}, \mathbb{C}^{k-i-1})$ .

**Proposition 4.76.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{K}[X_1, \dots, X_k]$ . Then*

- either  $1 \in \text{Ideal}(\mathcal{P}, \mathbb{K})$ ,
- or there exists a linear automorphism  $v: \mathbb{K}^k \rightarrow \mathbb{K}^k$  and a natural number  $k'$ ,  $0 \leq k' \leq k$ , such that the canonical projection  $\Pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k'}$  forgetting the last  $k - k'$  coordinates is a finite mapping from  $v(\text{Zer}(\mathcal{P}, \mathbb{C}^k))$  to  $\mathbb{C}^{k'}$  (the linear automorphism  $v$  being extended to  $\mathbb{C}^k$ ).

**Proof:** The proof is by induction on  $k$ .

When  $k = 1$ , consider the greatest common divisor  $Q$  of  $\mathcal{P}$ , which generates the ideal generated by  $\mathcal{P}$ . If  $Q = 0$  take  $k' = 1$ , and if  $Q \neq 0$ , take  $k' = 0$ .

If  $\text{deg}(Q) = 0$ , then  $Q$  is a non-zero constant and  $1 \in \text{Ideal}(\mathcal{P}, \mathbb{K})$ . If  $\text{deg}(Q) > 0$ , then  $\text{Zer}(\mathcal{P}, \mathbb{C}) = \text{Zer}(Q, \mathbb{C}^k)$  is non-empty and finite, so the projection to  $\{0\}$  is finite and the result holds in this case.

Suppose now that  $k > 1$  and that the theorem holds for  $k - 1$ .

If  $\text{Ideal}(\mathcal{P}, \mathbb{K}) = \{0\}$ , the theorem obviously holds by taking  $k' = 0$ .

If  $\text{Ideal}(\mathcal{P}, \mathbb{K}) \neq \{0\}$ , it follows from Lemma 4.73 that we can perform a linear change of variables  $w$  and assume that  $P_1(w(X))$  is quasi-monic with respect to  $X_k$ .

Let  $\mathcal{P}_w = \{P_1(w(X)), \dots, P_s(w(X))\}$ .

Applying the induction hypothesis to  $\text{Proj}_{X_k}(\mathcal{P}_w)$ ,

- either  $1 \in \text{Ideal}(\text{Proj}_{X_k}(\mathcal{P}_w, \mathbb{K}))$ ,
- or there exists a linear automorphism  $v': \mathbb{K}^{k-1} \rightarrow \mathbb{K}^{k-1}$  and a natural number  $k'$ ,  $0 \leq k' \leq k - 1$ , such that the canonical projection  $\Pi'$  from  $\mathbb{C}^{k-1}$  to  $\mathbb{C}^{k'}$  is a finite mapping from  $v'(\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}_w, \mathbb{C}^{k-1})))$  to  $\mathbb{C}^{k'}$ .

Since  $\text{Proj}_{X_k}(\mathcal{P}_w) \subset \text{Ideal}(\mathcal{P}, \mathbb{K})$ , it is clear that if  $1 \in \text{Ideal}(\text{Proj}_{X_k}(\mathcal{P}_w, \mathbb{K}))$ , then  $1 \in \text{Ideal}(\mathcal{P}_w, \mathbb{K})$ , which implies  $1 \in \text{Ideal}(\mathcal{P}, \mathbb{K})$ .

We now prove that if there exists a linear automorphism  $v': \mathbb{K}^{k-1} \rightarrow \mathbb{K}^{k-1}$  and a natural number  $k'$ ,  $0 \leq k' \leq k - 1$ , such that the canonical projection  $\Pi'$  from  $\mathbb{C}^{k-1}$  to  $\mathbb{C}^{k'}$  is a finite mapping from  $v'(\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}_w, \mathbb{C}^{k-1})))$  to  $\mathbb{C}^{k'}$ , there exists a linear automorphism  $v: \mathbb{K}^k \rightarrow \mathbb{K}^k$  such that the canonical projection  $\Pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k'}$  is a finite mapping from  $v(\text{Zer}(\mathcal{P}, \mathbb{C}^k))$  to  $\mathbb{C}^{k'}$ . By Proposition 4.75,  $w^{-1}(\text{Zer}(\mathcal{P}, \mathbb{C}^k)) = \text{Zer}(\mathcal{P}_w, \mathbb{C}^k)$  is finite on  $\text{Zer}(\text{Proj}_{X_k}(\mathcal{P}_w, \mathbb{C}^{k-1}))$ , so  $v = (v', \text{Id}) \circ w^{-1}$  is a linear automorphism from  $\mathbb{K}^k$  to  $\mathbb{K}^k$  such that the canonical projection  $\Pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k'}$  is a finite mapping from  $v(\text{Zer}(\mathcal{P}, \mathbb{C}^k))$  to  $\mathbb{C}^{k'}$ . □

We are now ready for the proof of Theorem 4.71 (Weak Hilbert's Nullstellensatz).

**Proof of Theorem 4.71:** The existence of  $A_1, \dots, A_s \in K[X_1, \dots, X_k]$  such that  $A_1 P_1 + \dots + A_s P_s = 1$  clearly implies  $\text{Zer}(\mathcal{P}, \mathbb{C}^k) = \emptyset$ .

On the other hand, by Proposition 4.76, if  $\text{Zer}(\mathcal{P}, \mathbb{C}^k) = \emptyset$ , there cannot exist a linear automorphism  $v: K^k \rightarrow K^k$  and a natural number  $k', 0 \leq k' \leq k$  such that the canonical projection  $\Pi$  from  $\mathbb{C}^k$  to  $\mathbb{C}^{k'}$  is a finite mapping from  $v(\text{Zer}(\mathcal{P}, \mathbb{C}^k))$  to  $\mathbb{C}^{k'}$ , since such a map must be surjective by definition.

So, using Proposition 4.76,  $1 \in \text{Ideal}(\mathcal{P}, K)$  which means that there exist  $A_1, \dots, A_s \in K[X_1, \dots, X_k]$  such that  $A_1 P_1 + \dots + A_s P_s = 1$ .  $\square$

Hilbert's Nullstellensatz is derived from the weak form of Hilbert's Nullstellensatz (Theorem 4.71) using what is commonly known as Rabinovitch's trick.

**Theorem 4.77. [Hilbert's Nullstellensatz]** *Let  $\mathcal{P}$  be a finite subset of  $K[X_1, \dots, X_k]$ . If a polynomial  $P$  with coefficients in  $K$  vanishes on  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , then  $P^n \in \text{Ideal}(\mathcal{P}, K)$  for some  $n$ .*

**Proof:** The set of polynomials  $\mathcal{P} \cup \{TP - 1\}$  in the variables  $X_1, \dots, X_k, T$  has no common zeros in  $\mathbb{C}^{k+1}$  so according to Theorem 4.71 if  $\mathcal{P} = \{P_1, \dots, P_s\}$ , there exist polynomials

$$A_1(X_1, \dots, X_k, T), \dots, A_s(X_1, \dots, X_k, T), A(X_1, \dots, X_k, T)$$

in  $K[X_1, \dots, X_k, T]$  such that  $1 = A_1 P_1 + \dots + A_s P_s + A(TP - 1)$ . Replacing everywhere  $T$  by  $1/P$  and clearing denominators by multiplying by an appropriate power of  $P$ , we see that a power of  $P$  is in the ideal  $\text{Ideal}(\mathcal{P}, K)$ .  $\square$

Another way of stating Hilbert's Nullstellensatz which follows immediately from the above is:

**Theorem 4.78.** *Let  $\mathcal{P}$  be a finite subset of  $K[X_1, \dots, X_k]$ . The radical of  $\text{Ideal}(\mathcal{P}, K)$  coincides with the set of polynomials in  $K[X_1, \dots, X_k]$  vanishing on  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  i.e.*

$$\sqrt{\text{Ideal}(\mathcal{P}, K)} = \{P \in K[X_1, \dots, X_k] \mid \forall x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k), P(x) = 0\}.$$

**Corollary 4.79. [Homogeneous Hilbert's Nullstellensatz]**

*Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset K[X_1, \dots, X_k]$  be a finite set of homogeneous polynomials with  $\deg(P_i) = d_i$ . If a homogeneous polynomial  $P \in K[X_1, \dots, X_k]$  of degree  $p$  vanishes on the common zeros of  $\mathcal{P}$  in  $\mathbb{C}^k$ , then there exists  $n \in \mathbb{N}$  and homogeneous polynomials  $H_1, \dots, H_s$  in  $K[X_1, \dots, X_k]$  of degrees  $c_1, \dots, c_s$  such that*

$$\begin{aligned} P^n &= H_1 P_1 + \dots + H_s P_s, \\ n p &= c_1 + d_1 = \dots = c_s + d_s. \end{aligned}$$

**Proof:** According to Hilbert's Nullstellensatz, there exist  $n \in \mathbb{N}$  and polynomials  $B_1, \dots, B_s$  in  $K[X_1, \dots, X_k]$  such that  $P^n = B_1 P_1 + \dots + B_s P_s$ .

Decompose  $B_i$  as  $H_i + C_i$  with  $H_i$  homogeneous of degree  $np - d_i$ , and notice that no monomial of  $C_i P_i$  has degree  $np$ . So  $P^n = H_1 P_1 + \dots + H_s P_s$ .  $\square$

*Remark 4.80.* Let us explain the statement claimed in the introduction of the chapter that our proof of Hilbert's Nullstellensatz is constructive.

Indeed, the method used for proving Theorem 4.71 (Weak Hilbert's Nullstellensatz) provides an algorithm for deciding, given a finite set  $\mathcal{P} \subset K[X_1, \dots, X_k]$ , whether  $\text{Zer}(\mathcal{P}, C^k)$  is empty and if it is empty, computes an algebraic identity certifying that 1 belongs to the ideal generated by  $\mathcal{P}$ .

The algorithm proceeds by eliminating variables one after the other. Given a finite family  $\mathcal{P} \neq \{0\}$  of polynomials in  $k$  variables, we check whether it contains a non-zero constant. If this is the case we conclude that  $\text{Zer}(\mathcal{P}, C^k)$  is empty and that 1 belongs to the ideal generated by  $\mathcal{P}$ . Otherwise, we perform a linear change of coordinates so that one of the polynomials of the family gets monic and replace the initial  $\mathcal{P}$  by this new family. Then we compute  $\text{Proj}_{X_k}(\mathcal{P})$ , which is a family of polynomials in  $k - 1$  variables together with algebraic identities expressing that the elements of  $\text{Proj}_{X_k}(\mathcal{P})$  belong to the ideal generated by  $\mathcal{P}$ . If  $\text{Proj}_{X_k}(\mathcal{P}) = \{0\}$  we conclude that  $\text{Zer}(\mathcal{P}, C^k)$  is not empty. If  $\text{Proj}_{X_k}(\mathcal{P}) \neq \{0\}$  we continue the process replacing  $k$  by  $k - 1$  and  $\mathcal{P}$  by  $\text{Proj}_{X_k}(\mathcal{P})$ . After eliminating all the variables, we certainly have that either the family of polynomials is  $\{0\}$  or it contains a non-zero constant, and we can conclude in both cases.  $\square$

Let us illustrate the algorithm outlined in the preceding remark by two examples.

*Example 4.81.* a) Let  $\mathcal{P} = \{X_2, X_2 + X_1, X_2 + 1\}$ . The first polynomial is monic with respect to  $X_2$ . We consider the resultant with respect to  $X_2$  of  $X_2$  and  $(U + 1)X_2 + UX_1 + 1$ , which is equal to  $UX_1 + 1$ . Thus  $\text{Proj}_{X_2}(\mathcal{P}) = \{X_1, 1\} \neq \{0\}$ , and contains a non-zero constant. Moreover  $1 = (X_2 + 1) - X_2$ . So we already proved that 1 belongs to the ideal generated by  $\mathcal{P}$  and  $\text{Zer}(\mathcal{P}, C^2) = \emptyset$ .

b) Let  $\mathcal{P} = \{X_2, X_2 + X_1, X_2 + 2X_1\}$ . The first polynomial is monic with respect to  $X_2$ . The resultant with respect to  $X_2$  of  $X_2$  and  $(U + 1)X_2 + (U + 2)X_1$  is equal to  $(U + 2)X_1$ . Thus  $\text{Proj}_{X_2}(\mathcal{P}) = \{X_1\} \neq \{0\}$ , contains a single element, and  $\text{Proj}_{X_1}(\text{Proj}_{X_2}(\mathcal{P})) = 0$ . Thus 1 does not belong to the ideal generated by  $\mathcal{P}$  and  $\text{Zer}(\mathcal{P}, C^2) \neq \emptyset$ . In fact,  $\text{Zer}(\mathcal{P}, C^2) = \{(0, 0)\}$ .  $\square$

Since the proof of Theorem 4.71 (Weak Hilbert's Nullstellensatz) is constructive, it is not surprising that it produces a bound on the degrees of the algebraic identity output. More precisely we have the following quantitative version of Hilbert's Nullstellensatz.

**Theorem 4.82. [Quantitative Hilbert's Nullstellensatz]**

Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{K}[X_1, \dots, X_k]$  be a finite set of less than  $d$  polynomials, of degrees bounded by  $d$ . If a polynomial  $P \in \mathbb{K}[X_1, \dots, X_k]$  of degree at most  $d$  in  $k$  variables vanishes on the common zeros of  $\mathcal{P}$  in  $\mathbb{C}^k$ , then there exists  $n \leq d(2d)^{2^{k+1}}$  and  $s$  polynomials  $B_1, \dots, B_s$  in  $\mathbb{K}[X_1, \dots, X_k]$  each of degree  $\leq d(2d)^{2^{k+1}}$  such that  $P^n = B_1 P_1 + \dots + B_s P_s$ .

**Proof:** The proof of the theorem follows from a close examination of the proofs of Proposition 4.76 and Theorem 4.77, using the notation of these proofs.

Suppose that  $\mathcal{P} = \{P_1, \dots, P_s\}$  has no common zeros in  $\mathbb{C}^k$ .

We consider first the case of 1 variable  $X$ . Since  $\text{Zer}(\mathcal{P}, \mathbb{C}) = \emptyset$ ,

$$\text{Res}_X(P_1, P_2 + UP_3 + \dots + U^{s-2}P_s) \in \mathbb{K}[U]$$

is not the zero polynomial, and we can find  $u \in \mathbb{K}$  such that

$$\text{Res}_X(P_1, P_2 + uP_3 + \dots + u^{s-2}P_s)$$

is a non-zero element of  $\mathbb{K}$ . By Proposition 1.9, there exist  $U$  and  $V$  of degree at most  $d$  such that

$$UP_1 + V(P_2 + uP_3 + \dots + u^{s-2}P_s) = 1,$$

which gives the identity

$$1 = UP_1 + VP_2 + uVP_3 + \dots + u^{s-2}VP_s$$

with  $\deg(U), \deg(V) \leq d$ .

Consider now the case of  $k$  variables. Since  $\text{Res}_{X_k}(P_1, R)$  is the determinant of the Sylvester matrix, which is of size at most  $2d$ , the degree of  $\text{Res}_{X_k}(P_1, R)$  with respect to  $X_1, \dots, X_{k-1}, U$  is at most  $2d^2$  (the entries of the Sylvester matrix are polynomials of degrees at most  $d$  in  $X_1, \dots, X_{k-1}, U$ ). So there are at most  $2d^2$  polynomials of degree  $2d^2$  in  $k-1$  variables to which the induction hypothesis is applied. Thus, the function  $g$  defined by

$$\begin{aligned} g(d, 1) &= d \\ g(d, k) &= g(2d^2, k-1) \end{aligned}$$

bounds the degree of the polynomials  $A_i$ . It is clear that  $(2d)^{2^k}$  is always bigger than  $g(d, k)$ .

For  $P \neq 1$ , using Rabinovitch's trick and apply the preceding bound to  $P_1, \dots, P_s, PT - 1$ , we get an identity

$$A_1 P_1 + \dots + A_s P_s + A(PT - 1) = 1,$$

with  $A_1, \dots, A_s, A$  of degree at most  $(2d)^{2^{k+1}}$ . Replacing  $T$  by  $1/P$  and removing denominators gives an identity

$$P^n = B_1 P_1 + \dots + B_s P_s$$

with  $n \leq (2d)^{2^{k+1}}$  and  $\deg(B_i) \leq d(2d)^{2^{k+1}}$ . □

The following corollary follows from Theorem 4.82 using the proof of Corollary 4.79.

**Corollary 4.83.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset K[X_1, \dots, X_k]$  be a finite set of less than  $d$  homogeneous polynomials of degrees  $d_i$  bounded by  $d$ . If a homogeneous polynomial  $P \in K[X_1, \dots, X_k]$  of degree  $p$  vanishes on the common zeros of  $\mathcal{P}$  in  $C^k$ , there exist  $n \in \mathbb{N}$ ,  $n \leq (2d)^{2^{k+1}}$ , and homogeneous polynomials  $H_1, \dots, H_s$  in  $K[X_1, \dots, X_k]$  of respective degrees  $c_1, \dots, c_s$  such that*

$$\begin{aligned} P^n &= H_1 P_1 + \dots + H_s P_s \\ n p &= c_1 + d_1 = \dots = c_s + d_s. \end{aligned}$$

*Remark 4.84.* Note that the double exponential degree bound in the number of variables obtained in Theorem 4.82 comes from the fact that the elimination of one variable between polynomials of degree  $d$  using resultant produces polynomials of degree  $d^2$ . A similar process occurs in Chapter 11 when we study cylindrical decomposition. □

## 4.5 Zero-dimensional Systems

Let  $\mathcal{P}$  be a finite subset of  $K[X_1, \dots, X_k]$ . The set of zeros of  $\mathcal{P}$  in  $C^k$

$$\text{Zer}(\mathcal{P}, C^k) = \{x \in C^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\}$$

is also called the set of **solutions** in  $C^k$  of the polynomial system of equations  $\mathcal{P} = 0$ . Abusing terminology, we speak of the solutions of a polynomial system  $\mathcal{P}$ . A system of polynomial equations  $\mathcal{P}$  is **zero-dimensional** if it has a finite number of solutions in  $C^k$ , i.e. if  $\text{Zer}(\mathcal{P}, C^k)$  is a non-empty finite set. We denote by  $\text{Ideal}(\mathcal{P}, K)$  the ideal generated by  $\mathcal{P}$  in  $K[X_1, \dots, X_k]$ .

We are going to prove that a system of polynomial equations

$$\mathcal{P} \subset K[X_1, \dots, X_k]$$

is zero-dimensional if and only if the  $K$ -vector space

$$A = K[X_1, \dots, X_k] / \text{Ideal}(\mathcal{P}, K)$$

is finite dimensional. The proof of this result relies on Hilbert's Nullstellensatz.

**Theorem 4.85.** *Let  $K$  be a field,  $C$  an algebraically closed field containing  $K$ , and  $\mathcal{P}$  a finite subset of  $K[X_1, \dots, X_k]$ .*

*The vector space  $A = K[X_1, \dots, X_k] / \text{Ideal}(\mathcal{P}, K)$  is of finite dimension  $m > 0$  if and only if  $\mathcal{P}$  is zero-dimensional.*

The number of elements of  $\text{Zer}(\mathcal{P}, C^k)$  is less than or equal to the dimension of  $A$  as a  $K$ -vector space.

Note that the fact that  $C$  is algebraically closed is essential in the statement, since otherwise there exist univariate polynomials of positive degree with no root.

The proof of the theorem will use the following definitions and results.

We consider the ideal  $\text{Ideal}(\mathcal{P}, C)$  generated by  $\mathcal{P}$  in  $C[X_1, \dots, X_k]$  and define  $\bar{A} = C[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, C)$ . Given  $x = (x_1, \dots, x_k) \in \text{Zer}(\mathcal{P}, C^k)$  and  $Q \in \bar{A}$ ,  $Q(x) \in C$  is well-defined, since two polynomials in  $C[X_1, \dots, X_k]$  having the same image in  $\bar{A}$  have the same value at  $x$ .

The following result makes precise the relationship between  $A$  and  $\bar{A}$ .

**Lemma 4.86.**  $A \subset \bar{A}$ .

**Proof:** If  $a$  and  $b$  are elements of  $K[X_1, \dots, X_k]$  equal modulo  $\text{Ideal}(\mathcal{P}, C)$ , then there exists for each  $P \in \mathcal{P}$  a polynomial  $A_P$  of some degree  $d_P$  in  $C[X_1, \dots, X_k]$  such that  $a - b = \sum A_P P$ . Since the various polynomials  $A_P P$  are linear combinations of a finite number of monomials, this identity can be seen as the fact that a system of linear equations with coefficients in  $K$  has a solution in  $C$  (the unknowns being the coefficients of the  $A_P$ ). We know from elementary linear algebra that this system of linear equations must then also have solutions in  $K$ , which means that there are  $C_P \in K[X_1, \dots, X_k]$  such that  $a - b = \sum_{P \in \mathcal{P}} C_P P$ . Thus  $a = b$  in  $A$ . This implies that the inclusion morphism  $A \subset \bar{A}$  is well-defined.  $\square$

We also have

**Lemma 4.87.** Let  $\mathcal{P}$  be a finite set of polynomials in  $K[X_1, \dots, X_k]$ . Then  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$  is a finite dimensional vector space of dimension  $m$  over  $K$  if and only if  $\bar{A} = C[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, C)$  is a finite dimensional vector space of dimension  $m$  over  $C$ .

**Proof:** Suppose that  $A$  has finite dimension  $m$  and consider any finite set of  $m' > m$  monomials in  $K[X_1, \dots, X_k]$ . Since the images in  $A$  of these monomials are linearly dependent in  $A$  over  $K$ , the images in  $\bar{A}$  of these monomials are linearly dependent in  $\bar{A}$  over  $C$ . Therefore the dimension of  $\bar{A}$  is finite and no greater than the dimension of  $A$ , since both  $A$  and  $\bar{A}$  are spanned by monomials.

For the other direction, if  $\bar{A}$  has finite dimension  $m$  then we consider any family  $B_1, \dots, B_m$  of  $m' > m$  elements in  $K[X_1, \dots, X_k]$  and denote by  $b_1, \dots, b_m$  their images in  $A$  and by  $b'_1, \dots, b'_{m'}$  their images in  $\bar{A}$ . Since  $b'_1, \dots, b'_{m'}$  are linearly dependent, there exist  $(\lambda_1, \dots, \lambda_{m'})$  in  $C^{m'}$  which are not all zero and for each  $P \in \mathcal{P}$  a polynomial  $A_P$  of some degree  $d_P$  in  $C[X_1, \dots, X_k]$  such that

$$\lambda_1 B_1 + \dots + \lambda_{m'} B_{m'} = \sum A_P P. \quad (4.5)$$

Since the various polynomials  $A_P P$  are linear combinations of a finite number of monomials, the identity 4.11 can be seen as the fact that a system of linear equations with coefficients in  $K$  has a solution in  $C$  (the unknowns being the  $\lambda_i$  and the coefficients of the  $A_P$ ). We know from elementary linear algebra that this system of linear equations must then also have solutions in  $K$  which means that there are  $\mu_i \in K$  not all zero and  $C_P \in K[X_1, \dots, X_k]$  such that

$$\mu_1 B_1 + \dots + \mu_{m'} B_{m'} = \sum C_P P.$$

Thus,  $b_1, \dots, b_{m'}$  are linearly dependent over  $K$  and hence the dimension of  $A$  is no greater than the dimension of  $\bar{A}$ . □

**Definition 4.88.** An element  $a$  of  $A$  is **separating for  $\mathcal{P}$**  if  $a$  has distinct values at distinct elements of  $\text{Zer}(\mathcal{P}, C^k)$ . □

Separating elements always exist when  $\mathcal{P}$  is zero-dimensional.

**Lemma 4.89.** *If  $\#\text{Zer}(\mathcal{P}, C^k) = n$ , then there exists  $i$ ,  $0 \leq i \leq (k - 1) \binom{n}{2}$ , such that*

$$a_i = X_1 + i X_2 + \dots + i^{k-1} X_k$$

*is separating.*

**Proof:** Let  $x = (x_1, \dots, x_k), y = (y_1, \dots, y_k)$  be two distinct points of  $\text{Zer}(\mathcal{P}, C^k)$  and let  $\ell(x, y)$  be the number of  $i$ ,  $0 \leq i \leq (k - 1) \binom{n}{2}$ , such that  $a_i(x) = a_i(y)$ . Since the polynomial

$$(x_1 - y_1) + (x_2 - y_2)T + \dots + (x_k - y_k)T^{k-1}$$

is not identically zero, it has no more than  $k - 1$  distinct roots. It follows that  $\ell(x, y) \leq k - 1$ . As the number of 2-element subsets of  $\text{Zer}(\mathcal{P}, C^k)$  is  $\binom{n}{2}$ , the lemma is proved. □

An important property of separating elements is the following lemma:

**Lemma 4.90.** *If  $a$  is separating and  $\text{Zer}(\mathcal{P}, C^k)$  has  $n$  elements, then  $1, a, \dots, a^{n-1}$  are linearly independent in  $A$ .*

**Proof:** Suppose that there exist  $c_i \in K$  such that  $\sum_{i=0}^{n-1} c_i a^i = 0$  in  $A$ , whence the polynomial  $c_0 + c_1 a + \dots + c_{n-1} a^{n-1}$  is in  $\text{Ideal}(\mathcal{P}, K)$ . Thus for all  $x \in \text{Zer}(\mathcal{P}, C^k)$ ,  $\sum_{i=0}^{n-1} c_i a^i(x) = 0$ . The univariate polynomial  $\sum_{i=0}^{n-1} c_i T^i = 0$  has  $n$  distinct roots and is therefore identically zero. □

**Proof of Theorem 4.85:** If  $A$  is a finite dimensional vector space of dimension  $N$  over  $K$ , then  $1, X_1, \dots, X_1^N$  are linearly dependent in  $A$ . Consequently, there is a polynomial  $p_1(X_1)$  of degree at most  $N$  in the ideal  $\text{Ideal}(\mathcal{P}, C)$ . It follows that the first coordinate of any  $x \in \text{Zer}(\mathcal{P}, C^k)$  is a root of  $p_1$ . Doing the same for all the variables, we see that  $\text{Zer}(\mathcal{P}, C^k)$  is a finite set.

Conversely, if  $\text{Zer}(\mathcal{P}, \mathbf{C}^k)$  is finite, take a polynomial  $p_1(X_1) \in \mathbf{C}[X_1]$  whose roots are the first coordinates of the elements of  $\text{Zer}(\mathcal{P}, \mathbf{C}^k)$ . According to Hilbert's Nullstellensatz (Theorem 4.77) a power of  $p_1$  belongs to the ideal  $\text{Ideal}(\mathcal{P}, \mathbf{C})$ . Doing the same for all variables, we see that for every  $i$ , there exists a polynomial of degree  $d_i$  in  $\mathbf{C}[X_i]$  in the ideal  $\text{Ideal}(\mathcal{P}, \mathbf{C})$ . It follows that  $\bar{A}$  has a basis consisting of monomials whose degree in  $X_i$  is less than  $d_i$ . Thus,  $\bar{A}$  is finite dimensional over  $\mathbf{C}$ . We conclude that  $A$  is finite dimensional over  $\mathbf{K}$  using Lemma 4.87.

Part b) of the theorem follows from Lemma 4.89 and Lemma 4.90.  $\square$

We now explain how the quotient ring  $\bar{A}$  splits into a finite number of local factors, one for each  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$ . These local factors are used to define the multiplicities of the solutions of the system of polynomial equations. In the case where all the multiplicities are one these local factors will be nothing but the field  $\mathbf{C}$  itself, and the projection onto the factor corresponding to an  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$  consists in sending an element of  $\bar{A}$  to its value at  $x$ .

We need a definition. A **local ring**  $B$  is a ring, such that for every  $a \in B$ , either  $a$  is invertible or  $1 + a$  is invertible. A field is always a local ring.

**Exercise 4.2.** A ring  $B$  is local if and only if has a unique maximal (proper) ideal which is the set of non-invertible elements.

Given a multiplicative subset  $S$  of a ring  $A$  (i.e. a subset of  $A$  closed under multiplication), we define an equivalence relation on ordered pairs  $(a, s)$  with  $a \in A$  and  $s \in S$  by  $(a, s) \sim (a', s')$  if and only if there exist  $t \in S$  such that  $t(a s' - a' s) = 0$ . The class of  $(a, s)$  is denoted  $a/s$ . The **ring of fractions**  $S^{-1}A$  is the set of classes  $a/s$  equipped with the following addition and multiplication

$$\begin{aligned} a/s + a'/s' &= (a s' + a' s)/(s s'), \\ (a/s)(a'/s') &= (a a')/(s s'). \end{aligned}$$

The **localization of  $\bar{A}$  at  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$**  is denoted  $\bar{A}_x$ . It is the ring of fractions associated to the multiplicative subset  $S_x$  consisting of elements of  $\bar{A}$  not vanishing at  $x$ . The ring  $\bar{A}_x$  is local: an element  $P/Q$  of  $\bar{A}_x$  is invertible if and only if  $P(x) \neq 0$ , and it is clear that either  $P/Q$  is invertible or  $1 + P/Q = (Q + P)/Q$  is invertible.

We will prove that the ring  $\bar{A}$  is isomorphic to the product of the various  $\bar{A}_x$  for  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$ . The proof relies on the following result.

**Proposition 4.91.** *If  $\text{Zer}(\mathcal{P}, \mathbf{C}^k)$  is finite, then, for every  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$ , there exists an element  $e_x$  of  $\bar{A}$  such that*

- $\sum_{x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)} e_x = 1$ ,
- $e_x e_y = 0$  for  $y \neq x$  with  $y, x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$ ,
- $e_x^2 = e_x$ ,
- $e_x(x) = 1$  for  $x \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$ ,
- $e_x(y) = 0$  for  $x, y \in \text{Zer}(\mathcal{P}, \mathbf{C}^k)$  and  $x \neq y$ .



**Proof:** We first prove that, for every  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , there exists an element  $s_x$  of  $\bar{A}$  with  $s_x(x) = 1$ ,  $s_x(y) = 0$  for every  $y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ ,  $y \neq x$ . Making if necessary a linear change of variables, we suppose that the variable  $X_1$  is separating. The classical Lagrange interpolation formula gives polynomials in  $X_1$  with the required properties. Namely, writing each  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$  as  $(x_1, \dots, x_k)$ , we set

$$s_x = \prod_{\substack{y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k) \\ y \neq x}} \frac{X_1 - y_1}{x_1 - y_1}.$$

Since  $s_x s_y$  vanishes at every common zero of  $\mathcal{P}$ , Hilbert's Nullstellensatz (Theorem 4.77) implies that there exists a power of each  $s_x$ , denoted  $t_x$ , such that  $t_x t_y = 0$  in  $\bar{A}$  for  $y \neq x$ , and  $t_x(x) = 1$ . The family of polynomials  $\mathcal{P} \cup \{t_x \mid x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)\}$  has no common zero so, according to Hilbert's Nullstellensatz, there exist polynomials  $r_x$  such that  $\sum t_x r_x = 1$  in  $\bar{A}$ . Take  $e_x = t_x r_x$ . It is now easy to verify the claimed properties.  $\square$

The element  $e_x$  is called the **idempotent associated to  $x$** . Since  $e_x^2 = e_x$ , the set  $e_x \bar{A}$  equipped with the restriction of the addition and multiplication of  $\bar{A}$  is a ring with identity (namely  $e_x$ ).

**Proposition 4.92.** *The ring  $e_x \bar{A}$  is isomorphic to the localization  $\bar{A}_x$  of  $\bar{A}$  at  $x$ .*

**Proof:** Note that if  $Q(x) \neq 0$ ,  $e_x Q$  is invertible in  $e_x \bar{A}$ . Indeed, we can decompose  $Q = Q(x)(1 + v)$  with  $v(x) = 0$ . Since  $\forall y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , we have  $v e_x(y) = 0$ ,  $(v e_x)^N = 0$  for some  $N \in \mathbb{N}$  by Hilbert's Nullstellensatz and thus  $e_x(1 + v)$  is invertible in  $e_x \bar{A}$ . Its inverse is

$$(1 - e_x v + \dots + (-v)^{N-1}) e_x,$$

and it follows that  $e_x Q$  is invertible as well.

So, denoting by  $(e_x Q)^{-1}$  the inverse of  $e_x Q$  in  $e_x \bar{A}$ , consider the mapping from  $\bar{A}_x$  to  $\bar{A}$  which sends  $P/Q$  to  $P(e_x Q)^{-1} = e_x P(e_x Q)^{-1}$ . It is easy to check that this is a ring homomorphism. Conversely, to  $P e_x$  is associated  $P/1$ , which is a ring homomorphism from  $e_x \bar{A}$  to  $\bar{A}_x$ .

To see that these two ring homomorphisms are inverses to each other, we need only prove that  $(P(e_x Q)^{-1})/1 = P/Q$  in  $\bar{A}_x$ . This is indeed the case since

$$P e_x ((e_x Q)(e_x Q)^{-1} - 1) = 0$$

and  $e_x(x) = 1$ .  $\square$

We now prove that  $\bar{A}$  is the product of the  $\bar{A}_x$ .

**Theorem 4.93.** *For each  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$  there exists an idempotent  $e_x$  such that  $e_x \bar{A} = \bar{A}_x$  and*

$$\bar{A} \cong \prod_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \bar{A}_x.$$

**Proof:** Since  $\sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} e_x = 1$ ,  $\bar{A} \cong \prod_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \bar{A}_x$ . The canonical surjection of  $\bar{A}$  onto  $\bar{A}_x$  coincides with multiplication by  $e_x$ .  $\square$

We denote by  $\mu(x)$  the dimension of  $\bar{A}_x$  as a  $\mathbb{C}$ -vector space. We call  $\mu(x)$  the **multiplicity** of the zero  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .

If the multiplicity of  $x$  is 1 we say that  $x$  is **simple**. Then  $\bar{A}_x = \mathbb{C}$  and the canonical surjection  $\bar{A}$  onto  $\bar{A}_x$  coincides with the homomorphism from  $\bar{A}$  to  $\mathbb{C}$  sending  $P$  to its value at  $x$ . Indeed, suppose that  $P(x) = 0$ . Then  $P e_x(y) = 0$  for every  $y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$  and hence by Hilbert's Nullstellensatz there exists  $N \in \mathbb{N}$  such that  $(P e_x)^N = 0$ . Since  $e_x$  is idempotent this implies that  $P^N e_x = 0$ , and thus  $P^N = 0$  in  $\bar{A}_x$  which is a field. Hence  $P = 0$  in  $\bar{A}_x$ .

When the system of polynomial equations  $\mathcal{P} = \{P_1, \dots, P_k\}$  is zero-dimensional, simple zeros coincide with non-singular zeros as we see now.

Let  $P_1, \dots, P_k$  be polynomials in  $\mathbb{C}[X_1, \dots, X_k]$ . A **non-singular zero** of

$$P_1(X_1, \dots, X_k), \dots, P_k(X_1, \dots, X_k)$$

is a  $k$ -tuple  $x = (x_1, \dots, x_k)$  of elements of  $\mathbb{C}^k$  such that

$$P_1(x_1, \dots, x_k) = \dots = P_k(x_1, \dots, x_k) = 0$$

and  $\det\left(\left[\frac{\partial P_i}{\partial X_j}(x)\right]\right) \neq 0$ .

**Proposition 4.94.** *Let  $\mathcal{P} = \{P_1, \dots, P_k\} \subset K[X_1, \dots, X_k]$  be a zero dimensional system and  $x$  a zero of  $\mathcal{P}$  in  $\mathbb{C}^k$ . Then the following are equivalent*

- $x$  is a non-singular zero of  $\mathcal{P}$ ,
- $x$  is simple, i.e. the multiplicity of  $x$  is 1 and  $\bar{A}_x = \mathbb{C}$ ,
- $M_x \subset \text{Ideal}(\mathcal{P}, \mathbb{C}) + (M_x)^2$ , denoting by  $M_x$  the ideal of elements of  $\mathbb{C}[X_1, \dots, X_k]$  vanishing at  $x$ .

**Proof:**  $a) \Rightarrow c)$  Using Taylor's formula at  $x$ ,

$$P_j = \sum_i \frac{\partial P_j}{\partial X_i}(x)(X_i - x_i) + B_j$$

with  $B_j \in (M_x)^2$ . So

$$\sum_i \frac{\partial P_j}{\partial X_i}(x)(X_i - x_i) \in \text{Ideal}(\mathcal{P}, \mathbb{K}) + (M_x)^2.$$

Since the matrix  $\left[\frac{\partial P_j}{\partial X_i}(x)\right]$  is invertible, for every  $i$

$$(X_i - x_i) \in \text{Ideal}(\mathcal{P}, \mathbb{K}) + (M_x)^2.$$

$c) \Rightarrow b)$  Since  $(X_i - x_i)e_x$  vanishes on  $\text{Zer}(\mathcal{P})$  for every  $i$ , and  $e_x^2 = e_x$ , according to Hilbert's Nullstellensatz, there exists  $N_i$  such that

$$(X_i - x_i)^{N_i} e_x \in \text{Ideal}(\mathcal{P}, \mathbb{K}).$$

So there exist  $N$  such that  $(M_x)^N \cdot e_x \subset \text{Ideal}(\mathcal{P}, \mathbb{K})$ . Using repeatedly  $M_x \subset \text{Ideal}(\mathcal{P}, \mathbb{K}) + (M_x)^2$ , we get

$$(M_x)^{N-1} \cdot e_x \subset \text{Ideal}(\mathcal{P}, \mathbb{K}), \dots, M_x \cdot e_x \subset \text{Ideal}(\mathcal{P}, \mathbb{K}).$$

This implies  $\overline{A}_x = \mathbb{C}$ .

$b) \Rightarrow a)$  If  $\overline{A}_x = \mathbb{C}$ , then for every  $i$ ,  $(X_i - x_i) e_x \in \text{Ideal}(\mathcal{P}, \mathbb{K})$ . Indeed  $(X_i - x_i) e_x = 0$  in  $\overline{A}_x = \mathbb{C}$  and  $(X_i - x_i) e_x e_y = 0$  in  $\overline{A}_y$  for  $y \neq x$  and  $y \in \text{Zer}(\mathcal{P}, \mathbb{C})$ . So, for every  $i$  there exist  $A_{i,j}$  such that

$$(X_i - x_i) e_x = \sum_j A_{i,j} P_j.$$

Differentiating with respect to  $X_i$  and  $X_\ell, \ell \neq i$  and evaluating at  $x$  we get

$$\begin{aligned} 1 &= \sum_j A_{i,j}(x) \frac{\partial P_j}{\partial X_i}(x), \\ 0 &= \sum_j A_{i,j}(x) \frac{\partial P_j}{\partial X_\ell}(x), \ell \neq i, \end{aligned}$$

so the matrix  $\left[ \frac{\partial P_j}{\partial X_i}(x) \right]$  is invertible. □

## 4.6 Multivariate Hermite's Quadratic Form

We consider a zero dimensional system  $\mathcal{P}$  and its set of solutions in  $\mathbb{C}^k$

$$\text{Zer}(\mathcal{P}, \mathbb{C}^k) = \{x \in \mathbb{C}^k \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\}.$$

We indicate first the relations between  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  and the eigenvalues of certain linear maps on the finite dimensional vector spaces

$$\begin{aligned} A &= \mathbb{K}[X_1, \dots, X_k] / \text{Ideal}(\mathcal{P}, \mathbb{K}) \quad \text{and} \\ \overline{A} &= \mathbb{C}[X_1, \dots, X_k] / \text{Ideal}(\mathcal{P}, \mathbb{C}). \end{aligned}$$

**Notation 4.95. [Multiplication map]** If  $f \in A$ , we denote by  $L_f: A \rightarrow A$  the linear map of multiplication by  $f$  defined by  $L_f(g) = fg$  for  $g \in A$ . Similarly, if  $f \in \overline{A}$ , we also denote by  $L_f: \overline{A} \rightarrow \overline{A}$  the linear map of multiplication by  $f$  defined by  $L_f(g) = fg$  for  $g \in \overline{A}$ . By Lemma 4.86,  $A \subset \overline{A}$ , so we denote also by  $L_f: \overline{A} \rightarrow \overline{A}$  the linear map of multiplication by  $f \in A$  defined by  $L_f(g) = fg$  for  $g \in \overline{A}$  and for  $f \in A$ . □

We denote as above by  $\overline{A}_x$  the localization at a zero  $x$  of  $\mathcal{P}$  and by  $\mu(x)$  its multiplicity. We denote by  $L_{f,x}$  the linear map of multiplication by  $f$  from  $\overline{A}_x$  to  $\overline{A}_x$  defined by  $L_{f,x}(P/Q) = fP/Q$ . Note that  $L_{f,x}$  is well-defined since if  $P_1/Q_1 = P/Q$ , then  $fP_1/Q_1 = fP/Q$ . Considering  $\overline{A}_x$  as a sub-vector space of  $\overline{A}$ ,  $L_{f,x}$  is the restriction of  $L_f$  to  $\overline{A}_x$ .

**Theorem 4.96.** *The eigenvalues of  $L_f$  are the  $f(x)$ , with multiplicity  $\mu(x)$ , for  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .*

**Proof:** As  $e_x(f - f(x))$  vanishes on  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , Hilbert's Nullstellensatz (Theorem 4.77) implies that there exists  $m \in \mathbb{N}$  such that  $(e_x(f - f(x)))^m = 0$  in  $\bar{A}$ , which means that  $L_{e_x(f - f(x))}$  is nilpotent and has a unique eigenvalue 0 with multiplicity  $\mu(x)$ . Thus  $L_{f,x}$  has only one eigenvalue  $f(x)$  with multiplicity  $\mu(x)$ . Using Theorem 4.93 completes the proof.  $\square$

It follows immediately:

**Theorem 4.97. [Stickelberger]** *For  $f \in \bar{A}$ , the linear map  $L_f$  has the following properties:*

*The trace of  $L_f$  is*

$$\text{Tr}(L_f) = \sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \mu(x) f(x). \quad (4.6)$$

*The determinant of  $L_f$  is*

$$\det(L_f) = \prod_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} f(x)^{\mu(x)}. \quad (4.7)$$

*The characteristic polynomial  $\chi(\mathcal{P}, f, T)$  of  $L_f$  is*

$$\chi(\mathcal{P}, f, T) = \prod_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} (T - f(x))^{\mu(x)}. \quad (4.8)$$

Note that the statement on the trace is a generalization of Proposition 4.54.

*Remark 4.98.* Note that if  $f \in A$ ,  $\text{Tr}(L_f)$  and  $\det(L_f)$  are in  $K$  and  $\chi(\mathcal{P}, f, T) \in K[T]$ . Moreover, if the multiplication table of  $A$  in the basis  $\mathcal{B}$  has entries in a ring  $D$  contained in  $K$  and  $f$  has coefficients in  $D$  in the basis  $\mathcal{B}$ ,  $\text{Tr}(L_f)$  and  $\det(L_f)$  are in  $D$  and  $\chi(\mathcal{P}, f, T) \in D[T]$ .  $\square$

A consequence of Theorem 4.97 (Stickelberger) is a multivariate generalization of the univariate Hermite's theorem seen earlier in this chapter (Theorem 4.57).

We first define the multivariate generalization of Hermite quadratic form. For every  $Q \in A$ , we define the **Hermite's bilinear map** as the bilinear map:

$$\begin{aligned} \text{her}(\mathcal{P}, Q): A \times A &\longrightarrow K \\ (f, g) &\longmapsto \text{Tr}(L_{fg}Q). \end{aligned}$$

The corresponding quadratic form associated to  $\text{her}(\mathcal{P}, Q)$  is called the **Hermite's quadratic form**,

$$\begin{aligned} \text{Her}(\mathcal{P}, Q): A &\longrightarrow K \\ f &\longmapsto \text{Tr}(L_{f^2}Q). \end{aligned}$$

When  $Q = 1$  we simply write  $\text{her}(\mathcal{P}) = \text{her}(\mathcal{P}, 1)$  and  $\text{Her}(\mathcal{P}) = \text{Her}(\mathcal{P}, 1)$ .

We shall write  $A_{\text{rad}}$  to denote the ring  $K[X_1, \dots, X_k] / \sqrt{\text{Ideal}(\mathcal{P}, K)}$ .

The next theorem gives the connection between the radical of  $\text{Ideal}(\mathcal{P}, K)$  and the radical of the quadratic form  $\text{Her}(\mathcal{P})$ :

$$\text{Rad}(\text{Her}(\mathcal{P})) = \{f \in A \mid \forall g \in A \text{ her}(\mathcal{P})(f, g) = 0\}.$$

**Theorem 4.99.**

$$\sqrt{\text{Ideal}(\mathcal{P}, K)} = \text{Rad}(\text{Her}(\mathcal{P})).$$

**Proof:** Let  $f$  be an element of  $\sqrt{\text{Ideal}(\mathcal{P}, K)}$ . Then  $f$  vanishes on every element of  $\text{Zer}(\mathcal{P}, C^k)$ . So, applying Corollary 4.97, we obtain the following equality for every  $g \in K[X_1, \dots, X_k]$ :

$$\text{her}(\mathcal{P})(f, g) = \text{Tr}(L_{fg}) = \sum_{x \in \text{Zer}(\mathcal{P}, C^k)} \mu(x) f(x)g(x) = 0.$$

Conversely, if  $f$  is an element such that  $\text{her}(\mathcal{P})(f, g) = 0$  for any  $g$  in  $A$ , then Corollary 4.97 gives:

$$\forall g \in A \text{ her}(\mathcal{P})(f, g) = \text{Tr}(L_{fg}) = \sum_{x \in \text{Zer}(\mathcal{P}, C^k)} \mu(x) f(x)g(x) = 0. \tag{4.9}$$

Let  $a$  be a separating element (Definition 4.88). If  $\text{Zer}(\mathcal{P}, C^k) = \{x_1, \dots, x_n\}$ , Equality (4.15) used with each of  $g = 1, \dots, a^{n-1}$  gives,

$$\begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ a(x_1)^{n-1} & \dots & a(x_n)^{n-1} \end{bmatrix} \cdot \begin{bmatrix} \mu(x_1) f(x_1) \\ \vdots \\ \mu(x_n) f(x_n) \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

so that  $f(x_1) = \dots = f(x_n) = 0$ , since  $a$  is separating and the matrix at the left hand side is a Vandermonde matrix, hence invertible. Using Hilbert's Nullstellensatz 4.71, we obtain  $f \in \sqrt{\text{Ideal}(\mathcal{P}, K)}$  as desired. □

The following result generalizes Hermite's Theorem (Theorem 4.57) and has a very similar proof.

We define the Tarski-query of  $Q$  for  $\mathcal{P}$  as

$$\text{TaQ}(Q, \mathcal{P}) = \sum_{x \in \text{Zer}(\mathcal{P}, R^k)} \text{sign}(Q(x))$$

**Theorem 4.100. [Multivariate Hermite]**

$$\begin{aligned} \text{Rank}(\text{Her}(\mathcal{P}, Q)) &= \#\{x \in \text{Zer}(\mathcal{P}, C^k) \mid Q(x) \neq 0\}, \\ \text{Sign}(\text{Her}(\mathcal{P}, Q)) &= \text{TaQ}(Q, \mathcal{P}). \end{aligned}$$

**Proof:** Consider a separating element  $a$ . The elements  $1, a, \dots, a^{n-1}$  are linearly independent in  $A$  by Lemma 4.90 and can be completed to a basis  $\omega_1 = 1, \omega_2 = a, \dots, \omega_n = a^{n-1}, \omega_{n+1}, \dots, \omega_N$  of the  $K$ -vector space  $A$ .

Corollary 4.97 provides the following expression for the quadratic form  $\text{Her}(\mathcal{P}, Q)$ : if  $f = \sum_{j=1}^N f_j \omega_j \in \mathbb{A}$

$$\text{Her}(\mathcal{P}, Q)(f) = \sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \mu(x) Q(x) \left( \sum_{j=1}^N f_j \omega_j(x) \right)^2.$$

Consequently, denoting  $\text{Zer}(\mathcal{P}, \mathbb{C}^k) = \{x_1, \dots, x_n\}$ ,  $\text{Her}(\mathcal{P}, Q)$  is the map

$$f \mapsto (f_1, \dots, f_N) \cdot \Gamma^t \cdot \Delta(\mu(x_1) Q(x_1), \dots, \mu(x_n) Q(x_n)) \cdot \Gamma \cdot (f_1, \dots, f_N)^t,$$

where

$$\Gamma = \begin{bmatrix} 1 & a(x_1) & \dots & a(x_1)^{n-1} & \omega_{n+1}(x_1) & \dots & \omega_N(x_1) \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ 1 & a(x_n) & \dots & a(x_n)^{n-1} & \omega_{n+1}(x_n) & \dots & \omega_N(x_n) \end{bmatrix}$$

and  $\Delta$  denotes a diagonal matrix with the indicated diagonal entries. Therefore it suffices to prove that the rank of  $\Gamma$  is equal to  $n$ . But  $a$  is separating and the principal minor of the matrix  $\Gamma$  is a Vandermonde determinant.

Given  $(f_1, \dots, f_N)$ , let  $f = \sum_{i=1}^N f_i \omega_i$ . According to Corollary 4.97, the quadratic form  $\text{Her}(\mathcal{P}, Q)$  is given in this basis by

$$\sum_{y \in \text{Zer}(\mathcal{P}, \mathbb{R}^k)} \mu(y) Q(y) \left( \sum_{i=1}^N f_i \omega_i(y) \right)^2 + \sum_{z \in \text{Zer}(\mathcal{P}, \mathbb{C}^k) \setminus \text{Zer}(\mathcal{P}, \mathbb{R}^k)} \mu(z) Q(z) \left( \sum_{i=1}^N f_i \omega_i(z) \right)^2$$

as a quadratic form in the variables  $f_i$ . We have already seen in the first part of the proof that the  $n$  rows of  $\Gamma$  are linearly independent over  $\mathbb{C}$ . Moreover, if  $z$  and  $\bar{z}$  are complex conjugate solutions of  $\mathcal{P}$ , with  $Q(z) \neq 0$ ,

$$\mu(z) Q(z) \left( \sum_{i=1}^N f_i \omega_i(z) \right)^2 + \mu(\bar{z}) Q(\bar{z}) \left( \sum_{i=1}^N f_i \omega_i(\bar{z}) \right)^2$$

is easily seen to be a difference of two squares of real linear forms. Indeed, writing  $\mu(z) Q(z) = (a(z) + i b(z))^2$ ,

$$(a(z) + i b(z)) \left( \sum_{i=1}^N f_i \omega_i(x) \right) = L_{1,z} + i L_{2,z},$$

with  $s_i(z)$  and  $t_i(z)$  the real and imaginary part of  $\omega_i(z)$ ,

$$L_{1,z} = \sum_{i=1}^N (a(z) s_i(z) - b(z) t_i(z)) f_i$$

$$L_{2,z} = \sum_{i=1}^N (a(z) t_i(z) + b(z) s_i(z)) f_i$$

are real linear forms in  $f_1, \dots, f_N$  with coefficients in  $\mathbb{R}$  so that

$$\mu(z)Q(z)\left(\sum_{i=1}^N f_i\omega_i(z)\right)^2 + \mu(\bar{z})Q(\bar{z})\left(\sum_{i=1}^N f_i\omega_i(\bar{z})\right)^2 = 2L_{1,z}^2 - 2L_{2,z}^2.$$

Moreover,  $L(y, f), L_1(z), L_2(z)$  ( $y \in \text{Zer}(\mathcal{P}, \mathbb{R}^k), z, \bar{z} \in \text{Zer}(\mathcal{P}, \mathbb{C}^k) \setminus \text{Zer}(\mathcal{P}, \mathbb{R}^k)$ ) are linearly independent linear forms.

So the signature of  $\text{Her}(\mathcal{P}, Q)$  is the signature of

$$\sum_{y \in \text{Zer}(\mathcal{P}, \mathbb{R}^k)} \mu(y)Q(y)L(y, f)^2.$$

Since the linear forms  $L(y, f)$ , are linearly independent the signature of  $\text{Her}(\mathcal{P}, Q)$  is  $\sum_{y \in \text{Zer}(\mathcal{P}, \mathbb{R}^k)} \text{sign}(Q(y)) = \text{TaQ}(Q, \mathcal{P})$ .  $\square$

### 4.7 Projective Space and a Weak Bézout's Theorem

Let  $\mathbb{R}$  be a real closed field and  $\mathbb{C} = \mathbb{R}[i]$ . The **complex projective space of dimension  $k$** ,  $\mathbb{P}_k(\mathbb{C})$ , is the set of lines of  $\mathbb{C}^{k+1}$  through the origin.

A  $(k + 1)$ -tuple  $x = (x_0, x_1, \dots, x_k) \neq (0, 0, \dots, 0)$  of elements of  $\mathbb{C}$  defines a line  $\bar{x}$  through the origin. This is denoted by  $\bar{x} = (x_0: x_1: \dots: x_k)$  and  $(x_0, x_1, \dots, x_k)$  are **homogeneous coordinates** of  $\bar{x}$ . Clearly,

$$(x_0: x_1: \dots: x_k) = (y_0: y_1: \dots: y_k)$$

if and only if there exists  $\lambda \neq 0$  in  $\mathbb{C}$  with  $x_i = \lambda y_i$ .

A polynomial  $P$  in  $\mathbb{C}[X_{1,0}, \dots, X_{1,k_1}, \dots, X_{m,0}, \dots, X_{m,k_m}]$  is **multi-homogeneous** of multidegree  $d_1, \dots, d_m$  if it is homogeneous of degree  $d_i$  in the block of variables  $X_{i,0}, \dots, X_{i,k_i}$  for every  $i \leq m$ .

For example  $T(X^2 + Y^2)$  is homogeneous of degree 1 in the variable  $T$  and homogeneous of degree 2 in the variables  $\{X, Y\}$ .

If  $P$  is multi-homogeneous of multidegree  $d_1, \dots, d_m$ , a **zero** of  $P$  in  $\mathbb{P}_{k_1}(\mathbb{C}) \times \dots \times \mathbb{P}_{k_m}(\mathbb{C})$  is a point

$$x = (\bar{x}_1, \dots, \bar{x}_m) \in \mathbb{P}_{k_1}(\mathbb{C}) \times \dots \times \mathbb{P}_{k_m}(\mathbb{C})$$

such that  $P(x_1, \dots, x_m) = 0$ , and this property denoted by  $P(x) = 0$  depends only on  $x$  and not on the choice of the homogeneous coordinates. An **algebraic set** of  $\mathbb{P}_{k_1}(\mathbb{C}) \times \dots \times \mathbb{P}_{k_m}(\mathbb{C})$  is a set of the form

$$\text{Zer}(\mathcal{P}, \prod_{i=1}^m \mathbb{P}_{k_i}(\mathbb{C})) = \{x \in \prod_{i=1}^m \mathbb{P}_{k_i}(\mathbb{C}) \mid \bigwedge_{P \in \mathcal{P}} P(x) = 0\},$$

where  $\mathcal{P}$  is a finite set of multi-homogeneous polynomials in

$$\mathbb{C}[X_1, \dots, X_m] = \mathbb{C}[X_{1,0}, \dots, X_{1,k_1}, \dots, X_{m,0}, \dots, X_{m,k_m}].$$

**Lemma 4.101.** *An algebraic subset of  $\mathbb{P}_1(\mathbb{C})$  is either  $\mathbb{P}_1(\mathbb{C})$  or a finite set of points.*

**Proof:** Let  $\mathcal{P} = \{P_1, \dots, P_s\}$  with  $P_i$  homogeneous of degree  $d_i$ . If all the  $P_i$  are zero,  $\text{Zer}(\mathcal{P}, \mathbb{P}_1(\mathbb{C})) = \mathbb{P}_1(\mathbb{C})$ . Otherwise,  $\text{Zer}(\mathcal{P}, \mathbb{P}_1(\mathbb{C}))$  contains  $(0:1)$  if and only if

$$P_1(0, X_1) = \dots = P_s(0, X_1) = 0.$$

The other elements of  $\text{Zer}(\mathcal{P}, \mathbb{P}_1(\mathbb{C}))$  are the points of the form  $(1:x_1)$ , where  $x_1$  is a solution of

$$P_1(1, X_1) = \dots = P_s(1, X_1) = 0,$$

which is a finite number of points since the  $P_i(1, X_1)$  are not all zero.  $\square$

**Theorem 4.102.** *If  $V \subset \mathbb{P}_{k_1}(\mathbb{C}) \times \mathbb{P}_{k_2}(\mathbb{C})$  is algebraic, its projection on  $\mathbb{P}_{k_2}(\mathbb{C})$  is algebraic.*

**Proof:** We first introduce some notation. With  $X = (X_0, \dots, X_k)$ , we denote the set of homogeneous polynomials of degree  $d$  in  $X$  by  $C[X]_d$ . Let  $\mathcal{P} = \{P_1, \dots, P_s\}$  be a finite set of homogeneous polynomials with  $P_i$  of degree  $d_i$  in  $X$ . For  $d \geq d_i$ , let  $M_d(\mathcal{P})$  be the mapping

$$C[X]_{d-d_1} \times \dots \times C[X]_{d-d_s} \rightarrow C[X]_d$$

sending  $(H_1, \dots, H_s)$  to  $H_1 P_1 + \dots + H_s P_s$ . Identifying a homogeneous polynomial with its vector of coefficients in the basis of monomials,  $M_d(\mathcal{P})$  defines a matrix  $\mathcal{M}_d(\mathcal{P})$ .

The projection of  $\text{Zer}(\mathcal{P}, \mathbb{P}_{k_1}(\mathbb{C}) \times \mathbb{P}_{k_2}(\mathbb{C}))$  on  $\mathbb{P}_{k_2}(\mathbb{C})$  is

$$\pi(\text{Zer}(\mathcal{P}, \mathbb{P}_{k_1}(\mathbb{C}) \times \mathbb{P}_{k_2}(\mathbb{C}))) = \{\bar{y} \in \mathbb{P}_{k_2}(\mathbb{C}) \mid \exists \bar{x} \in \mathbb{P}_{k_1}(\mathbb{C}) \bigwedge_{P \in \mathcal{P}} P(x, y) = 0\}$$

Consider  $\bar{y} \notin \pi(V)$ , i.e.  $\bar{y} \in \mathbb{P}_{k_2}(\mathbb{C})$  and such that

$$\{\bar{x} \in \mathbb{P}_{k_1}(\mathbb{C}) \mid \bigwedge_{P \in \mathcal{P}} P(x, y) = 0\} = \emptyset.$$

Then

$$\{x \in C^{k_1+1} \mid \bigwedge_{P \in \mathcal{P}} P(x, y) = 0\} = \{0\}.$$

According to Corollary 4.83, there exists for every  $i = 0, \dots, k_i$  an integer  $n_i \leq (2d)^{2^{k_1+2}}$  and polynomials  $A_{i,j} \in C[X]_{n_i-d_i}$  such that

$$X_i^{n_i} = A_{i,1}(X) P_1(X, y) + \dots + A_{i,s}(X) P_s(X, y).$$

Since any monomial of degree  $N = \sum_{i=0}^{k_1} (2d)^{2^{k_1+2}}$  is a multiple of  $X_i^{n_i}$  for some  $1 \leq i \leq k_1$ , for every polynomial  $P \in C[X]_N$  there exist polynomials  $H_1, \dots, H_s$  with  $H_i \in C[X]_{N-d_i}$  such that

$$P = H_1(X) P_1(X, y) + \dots + H_s(X) P_s(X, y).$$



We have proved that  $\bar{y} \notin \pi(V)$  if and only if  $M_N(\{P_1(X, y), \dots, P_s(X, y)\})$  is surjective. This can be expressed by a finite disjunction of conditions  $M_i(y) \neq 0$ , where the  $M_i(Y)$  are the maximal minors extracted from the matrix  $\mathcal{M}_N(\{P_1(X, Y), \dots, P_s(X, Y)\})$  in which  $Y = (Y_0, \dots, Y_{k_2})$  appear as variables. Hence,  $\pi(V) = \{\bar{y} \mid \bigwedge M_i(y) = 0\}$ , which is an algebraic set of  $\mathbb{P}_{k_2}(\mathbb{C})$ .  $\square$

The remainder of the chapter is devoted to proving a weak version of Bézout's theorem, estimating the number of non-singular projective zeros of a polynomial system of equations. The proof of this theorem is quite simple. The basic idea is that we look at a polynomial system which has exactly the maximum number of such zeros and move continuously from this polynomial system to the one under consideration in such a way that the number of non-singular projective zeros cannot increase. In order to carry out this simple idea, we define projective zeros and elaborate a little on the geometry of  $\mathbb{P}_k(\mathbb{C})$ .

If  $P_1, \dots, P_k$  are homogeneous polynomials in  $\mathbb{C}[X_0, \dots, X_k]$ , we say that  $x = (x_0: x_1: \dots: x_k) \in \mathbb{P}_k(\mathbb{C})$  is a **non-singular projective zero** of  $P_1, \dots, P_k$  if  $P_i(x) = 0$  for  $i = 1, \dots, k$  and

$$\text{rank}\left(\left[\frac{\partial P_i}{\partial X_j}(x)\right]\right) = k,$$

for  $i = 1, \dots, k$ ,  $j = 0, \dots, k$ . Note that  $(x_1, \dots, x_k)$  is a non-singular zero of

$$P_1(1, X_1, \dots, X_k), \dots, P_k(1, X_1, \dots, X_k)$$

if and only if  $(1: x_1: \dots: x_k)$  is a non-singular projective zero of

$$P_1, \dots, P_k.$$

For  $i = 0, \dots, k$ , let  $\phi_i$  be the map from  $\mathbb{C}^k$  to  $\mathbb{P}_k(\mathbb{C})$  which maps  $(x_1, \dots, x_k)$  to  $(x_1: \dots: x_{i-1}: 1: x_i: \dots: x_k)$ , and let  $\mathcal{U}_i = \phi_i(\mathbb{C}^k)$ . Note that

$$\begin{aligned} \mathcal{U}_i &= \{\bar{x} \in \mathbb{P}_k(\mathbb{C}) \mid x_i \neq 0\}, \\ \phi_i^{-1}(x_0: x_{i-1}: x_i: x_{i+1}: \dots: x_k) &= \left(\frac{x_0}{x_i}, \dots, \frac{x_{i-1}}{x_i}, \frac{x_{i+1}}{x_i}, \dots, \frac{x_k}{x_i}\right). \end{aligned}$$

It is clear that  $\cup_{i=0, \dots, k} \mathcal{U}_i = \mathbb{P}_k(\mathbb{C})$ . It is also clear that  $\phi_i^{-1}(\mathcal{U}_i \cap \mathcal{U}_j)$  is a semi-algebraic open subset of  $\mathbb{C}^k = \mathbb{R}^{2k}$  and that  $\phi_j^{-1} \circ \phi_i$  is a semi-algebraic bijection from  $\phi_i^{-1}(\mathcal{U}_i \cap \mathcal{U}_j)$  to  $\phi_j^{-1}(\mathcal{U}_i \cap \mathcal{U}_j)$ .

We define the euclidean topology and semi-algebraic sets of  $\mathbb{P}_k(\mathbb{C})$  as follows:

- a subset  $U$  of  $\mathbb{P}_k(\mathbb{C})$  is **open** in the euclidean topology if only if for every  $i = 0, \dots, k$ ,  $\phi_i^{-1}(U \cap \mathcal{U}_i)$  is an open set in the euclidean topology of  $\mathbb{C}^k = \mathbb{R}^{2k}$ ,
- a subset  $S$  of  $\mathbb{P}_k(\mathbb{C})$  is **semi-algebraic** if only if for every  $i = 0, \dots, k$ ,  $\phi_i^{-1}(S \cap \mathcal{U}_i)$  is semi-algebraic in  $\mathbb{C}^k = \mathbb{R}^{2k}$ .

Note that the  $\mathcal{U}_i$  are semi-algebraic open subsets of  $\mathbb{P}_k(\mathbb{C})$ .

Similarly, it is easy to define the semi-algebraic sets of  $\mathbb{P}_k(\mathbb{C}) \times \mathbb{P}_\ell(\mathbb{C})$ . A semi-algebraic mapping from  $\mathbb{P}_k(\mathbb{C})$  to  $\mathbb{P}_\ell(\mathbb{C})$  is a mapping whose graph is semi-algebraic.

Since every point of  $\mathbb{P}_k(\mathbb{C})$  has a neighborhood that is contained in some  $\mathcal{U}_i$ , the local properties of  $\mathbb{P}_k(\mathbb{C})$  are the same as the local properties of  $\mathbb{C}^k = \mathbb{R}^{2k}$ . In particular the notion of differentiability and the classes  $\mathcal{S}^m$  and  $\mathcal{S}^\infty$  can be defined in a similar way and the corresponding implicit function theorem remains valid.

**Theorem 4.103. [Projective Implicit Function Theorem]**

Let  $(x^0, y^0) \in \mathbb{P}_k(\mathbb{C}) \times \mathbb{P}_\ell(\mathbb{C})$ , and let  $f_1, \dots, f_\ell$  be semi-algebraic functions of class  $\mathcal{S}^m$  on an open neighborhood of  $(\bar{x}^0, \bar{y}^0)$  such that  $f_j(x^0, y^0) = 0$  for  $j = 1, \dots, \ell$  and the Jacobian matrix

$$\left[ \frac{\partial f_j}{\partial y_i}(\bar{x}^0, \bar{y}^0) \right]$$

is invertible. Then there exists a semi-algebraic open neighborhood  $U$  (resp.  $V$ ) of  $\bar{x}^0$  (resp.  $\bar{y}^0$ ) in  $\mathbb{P}_k(\mathbb{C})$  (resp.  $\mathbb{P}_\ell(\mathbb{C})$ ) and a function  $\varphi \in \mathcal{S}^m(U, V)$  such that  $\varphi(\bar{x}^0) = \bar{y}^0$  and such that for every  $(\bar{x}, \bar{y}) \in U \times V$ , we have

$$f_1(\bar{x}, \bar{y}) = \dots = f_\ell(\bar{x}, \bar{y}) = 0 \Leftrightarrow \bar{y} = \varphi(\bar{x}).$$

Our final observation is the following lemma showing that the complement of a finite subset of  $\mathbb{P}_1(\mathbb{C})$  is semi-algebraically path connected.

If  $S$  is a semi-algebraic subset of  $\mathbb{P}_k(\mathbb{C})$ , we say that  $S$  is **semi-algebraically path connected** if for every  $x$  and  $y$  in  $S$ , there exists a continuous path from  $x$  to  $y$ , i.e. a continuous mapping  $\gamma$  from  $[0, 1]$  to  $S$  such that  $\gamma(0) = x, \gamma(1) = y$  and the graph of  $\gamma$  is semi-algebraic.

**Lemma 4.104.** *If  $\Delta$  is a finite subset of  $\mathbb{P}_1(\mathbb{C})$ , then  $\mathbb{P}_1(\mathbb{C}) \setminus \Delta$  is semi-algebraically path connected.*

**Proof:** If  $x$  and  $y$  both belong to  $\mathcal{U}_0$  (resp.  $\mathcal{U}_1$ ), it is clear that there is a semi-algebraic continuous path from  $\phi_0^{-1}(x)$  to  $\phi_0^{-1}(y)$  (resp.  $\phi_1^{-1}(x)$  to  $\phi_1^{-1}(y)$ ) avoiding  $\phi_0^{-1}(\Delta \cap \mathcal{U}_0)$  (resp.  $\phi_1^{-1}(\Delta \cap \mathcal{U}_0)$ ). If  $x \in \mathcal{U}_0, y \in \mathcal{U}_1$ , take  $z \in (\mathbb{P}_1(\mathbb{C}) \setminus \Delta) \cap \mathcal{U}_0 \cap \mathcal{U}_1$  and connect  $x$  to  $z$  and then  $z$  to  $y$  outside  $\Delta$  by semi-algebraic and continuous paths.  $\square$

**Proposition 4.105.** *Let  $P_1, \dots, P_k \in \mathbb{C}[X_0, \dots, X_k]$  be homogeneous polynomials of degrees  $d_1, \dots, d_k$ . Then the number of non-singular projective zeros of  $P_1, \dots, P_k$  is at most  $d_1 \cdots d_k$ .*

**Proof:** For  $i = 1, \dots, k$ , let

$$H_{i,\lambda,\mu}(X_0, \dots, X_k) = \lambda P_i + \mu (X_i - X_0) (X_i - 2X_0) \cdots (X_i - d_i X_0), \text{ for } (\lambda, \mu) \in \mathbb{C}^2 \setminus \{0\}.$$

We denote by  $\mathcal{S}_{(\lambda:\mu)}$  the polynomial system

$$H_{1,\lambda,\mu}, \dots, H_{k,\lambda,\mu}.$$

Note that the polynomial system  $\mathcal{S}_{(0:1)}$  has  $d_1 \cdots d_k$  non-singular projective zeros and  $\mathcal{S}_{(1:0)}$  is  $P_1, \dots, P_k$ . The subset of  $(x, (\lambda:\mu)) \in \mathbb{P}_k(\mathbb{C}) \times \mathbb{P}_1(\mathbb{C})$  such that  $x$  is a singular projective zero of the polynomial system  $\mathcal{S}_{(\lambda:\mu)}$  is clearly algebraic. Therefore, according to Theorem 4.102, its projection  $\Delta$  on  $\mathbb{P}_1(\mathbb{C})$  is an algebraic subset of  $\mathbb{P}_1(\mathbb{C})$ . Since  $(0:1) \notin \Delta$ , the set  $\Delta$  consists of finitely many points, using Lemma 4.101. Since  $\mathbb{P}_1(\mathbb{C}) \setminus \Delta$  is semi-algebraically connected, there is a semi-algebraic continuous path  $\gamma: [0, 1] \subset \mathbb{R} \rightarrow \mathbb{P}_1(\mathbb{C})$  such that  $\gamma(0) = (1:0)$ ,  $\gamma(1) = (0:1)$ , and  $\gamma((0, 1]) \subset \mathbb{P}_1(\mathbb{C}) \setminus \Delta$ . Note that  $(\lambda:\mu) \in \mathbb{P}_1(\mathbb{C}) \setminus \Delta$  if and only if all projective zeros of  $\mathcal{S}_{(\lambda:\mu)}$  are non-singular. By the implicit function theorem, for every non-singular projective zero  $x$  of  $\mathcal{S}_{(1:0)}$ , there exists a continuous path  $\sigma_x: [0, 1] \rightarrow \mathbb{P}_k(\mathbb{C})$  such that  $\sigma_x(0) = x$  and, for every  $t \in (0, 1]$ ,  $\sigma_x(t)$  is a non-singular projective zero of  $\mathcal{S}_{\gamma(t)}$ . Moreover, if  $y$  is another non-singular projective zero of  $\mathcal{S}_{(1:0)}$ , then  $\sigma_x(t) \neq \sigma_y(t)$  for every  $t \in [0, 1]$ . From this we conclude that the number of non-singular projective zeros of  $\mathcal{S}_{(1:0)}$ :  $P_1 = \dots = P_k = 0$  is less than or equal to the number of projective zeros of  $\mathcal{S}_{(0:1)}$ , which is  $d_1 \cdots d_k$ . □

**Theorem 4.106. [Weak Bézout]** *Let  $P_1, \dots, P_k \in \mathbb{C}[X_1, \dots, X_k]$  be polynomials of degrees  $d_1, \dots, d_k$ . Then the number of non-singular zeros of  $P_1, \dots, P_k$  is at most  $d_1 \cdots d_k$ .*

**Proof:** Define

$$P_i^h = X_0^{d_i} P_i \left( \frac{X_1}{X_0}, \dots, \frac{X_k}{X_0} \right), i = 1, \dots, k,$$

and apply Proposition 4.105. The claim follows, noticing that any non-singular zero of  $P_1, \dots, P_k$  is a non-singular projective zero of  $P_1^h, \dots, P_k^h$ . □

## 4.8 Bibliographical Notes

Resultants were introduced by Euler [56] and Bézout [24] and have been studied by many authors, particularly Sylvester [153]. Subresultant coefficients are discussed already in Gordan’s textbook [74].

The use of quadratic forms for real root counting, in the univariate and multivariate case, is due to Hermite [89].

Hilbert’s Nullstellensatz appears in [91] and a constructive proof giving doubly exponential degrees can be found in [88]. Much better degree bounds have been proved more recently, and are not included in our book [31, 35, 97].

---

## Decomposition of Semi-Algebraic Sets

In this chapter, we decompose semi-algebraic sets in various ways and study several consequences of these decompositions. In Section 5.1 we introduce the cylindrical decomposition which is a key technique for studying the geometry of semi-algebraic sets. In Section 5.2 we use the cylindrical decomposition to define and study the semi-algebraically connected components of a semi-algebraic set. In Section 5.3 we define the dimension of a semi-algebraic set and obtain some basic properties of dimension. In Section 5.4 we get a semi-algebraic description of the partition induced by the cylindrical decomposition using Thom's lemma. In Section 5.5 we decompose semi-algebraic sets into smooth manifolds, called strata, generalizing Thom's lemma in the multivariate case. In Section 5.6 we define simplicial complexes, and establish the existence of a triangulation for a closed and bounded semi-algebraic set in Section 5.7. This triangulation result is used in Section 5.8 to prove Hardt's triviality theorem which has several important consequences, notably among them the finiteness of topological types of algebraic sets defined by polynomials of fixed degrees. We conclude the chapter with a semi-algebraic version of Sard's theorem in Section 5.9.

### 5.1 Cylindrical Decomposition

We first define what is a cylindrical decomposition: a decomposition of  $\mathbb{R}^k$  into a finite number of semi-algebraically connected semi-algebraic sets having a specific structure with respect to projections.

**Definition 5.1.** A **cylindrical decomposition** of  $\mathbb{R}^k$  is a sequence  $\mathcal{S}_1, \dots, \mathcal{S}_k$  where, for each  $1 \leq i \leq k$ ,  $\mathcal{S}_i$  is a finite partition of  $\mathbb{R}^i$  into semi-algebraic subsets, called the **cells of level  $i$** , which satisfy the following properties:

- Each cell  $S \in \mathcal{S}_1$  is either a point or an open interval.
- For every  $1 \leq i < k$  and every  $S \in \mathcal{S}_i$ , there are finitely many continuous semi-algebraic functions

$$\xi_{S,1} < \dots < \xi_{S,\ell_S}: S \longrightarrow \mathbb{R}$$

such that the cylinder  $S \times \mathbb{R} \subset \mathbb{R}^{i+1}$  is the disjoint union of cells of  $\mathcal{S}_{i+1}$  which are:

- either the graph of one of the functions  $\xi_{S,j}$ , for  $j = 1, \dots, \ell_S$ :

$$\{(x', x_{j+1}) \in S \times \mathbb{R} \mid x_{j+1} = \xi_{S,j}(x')\},$$

- or a band of the cylinder bounded from below and from above by the graphs of the functions  $\xi_{S,j}$  and  $\xi_{S,j+1}$ , for  $j = 0, \dots, \ell_S$ , where we take  $\xi_{S,0} = -\infty$  and  $\xi_{S,\ell_S+1} = +\infty$ :

$$\{(x', x_{j+1}) \in S \times \mathbb{R} \mid \xi_{S,j}(x') < x_{j+1} < \xi_{S,j+1}(x')\}. \quad \square$$

*Remark 5.2.* Denoting by  $\pi_\ell$  the canonical projection of  $\mathbb{R}^k$  to  $\mathbb{R}^\ell$ , it follows immediately from the definition that for every cell  $T$  of  $\mathcal{S}_i$ ,  $i \geq \ell$ ,  $S = \pi_\ell(T)$  is a cell of  $\mathcal{S}_\ell$ . We say that the cell  $T$  lies above the cell  $S$ . It is also clear that if  $S$  is a cell of  $\mathcal{S}_i$ , denoting by  $T_1, \dots, T_m$  the cells of  $\mathcal{S}_{i+1}$  lying above  $S$ , we have  $S \times \mathbb{R} = \bigcup_{j=1}^m T_j$ .  $\square$

**Proposition 5.3.** *Every cell of a cylindrical decomposition is semi-algebraically homeomorphic to an open  $i$ -cube  $(0, 1)^i$  (by convention,  $(0, 1)^0$  is a point) and is semi-algebraically connected.*

**Proof:** We prove the proposition for the cells of  $\mathcal{S}_i$  by induction on  $i$ .

If  $i = 0$ , the cells are clearly either points or open intervals and the claim holds.

Observe that if  $S$  is a cell of  $\mathcal{S}_i$ , the graph of  $\xi_{S,j}$  is semi-algebraically homeomorphic to  $S$  and every band

$$\{(x', x_{j+1}) \in S \times \mathbb{R} \mid \xi_{S,j}(x') < x_{j+1} < \xi_{S,j+1}(x')\}$$

is semi-algebraically homeomorphic to  $S \times (0, 1)$ . In the case of the graph of  $\xi_{S,j}$ , the homeomorphism simply maps  $x' \in S$  to  $(x', \xi_{S,j}(x'))$ .

For  $\{(x', x_{j+1}) \in S \times \mathbb{R} \mid \xi_{S,j}(x') < x_{j+1} < \xi_{S,j+1}(x')\}$ ,  $0 < j < \ell_S$  we map  $(x', t) \in S \times (0, 1)$  to  $(x', (1-t)\xi_{S,j}(x') + t\xi_{S,j+1}(x'))$ .

In the special case  $j = 0$ ,  $j = \ell_S$ , we take

$$\begin{aligned} & \left( x', \frac{t-1}{t} + \xi_{S,j}(x') \right) \quad \text{if } j = 0, \ell_S \neq 0, \\ & \left( x', \frac{t}{1-t} + \xi_{S,\ell_S}(x') \right) \quad \text{if } j = \ell_S \neq 0, \\ & \left( x', -\frac{1}{t} + \frac{1}{1-t} \right) \quad \text{if } j = \ell_S = 0. \end{aligned}$$

These mappings are clearly bijective, bicontinuous and semi-algebraic, noting that the mappings sending  $t \in (0, 1)$  to

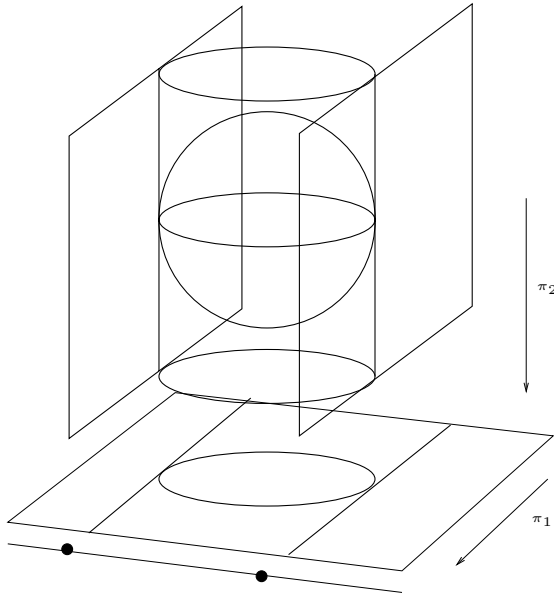
$$\frac{t-1}{t} + a, \frac{t}{1-t} + a, -\frac{1}{t} + \frac{1}{1-t},$$

are semi-algebraic bijections from  $(0, 1)$  to  $(-\infty, a)$ ,  $(a, +\infty)$ ,  $(-\infty, +\infty)$ .  $\square$

A **cylindrical decomposition adapted to a finite family of semi-algebraic sets**  $T_1, \dots, T_\ell$  is a cylindrical decomposition of  $\mathbb{R}^k$  such that every  $T_i$  is a union of cells.

*Example 5.4.* We illustrate this definition by presenting a cylindrical decomposition of  $\mathbb{R}^3$  adapted to the unit sphere.

Note that the projection of the sphere on the  $X_1, X_2$  plane is the unit disk. The intersection of the sphere and the cylinder above the open unit disk consists of two hemispheres. The intersection of the sphere and the cylinder above the unit circle consists of a circle. The intersection of the sphere and the cylinder above the complement of the unit disk is empty. Note also that the projection of the unit circle on the line is the interval  $[-1, 1]$ .



**Fig. 5.1.** A cylindrical decomposition adapted to the sphere in  $\mathbb{R}^3$

The decomposition of  $\mathbb{R}$  consists of five cells of level 1 corresponding to the points  $-1$  and  $1$  and the three intervals they define.

$$\begin{cases} S_1 = (-\infty, -1) \\ S_2 = \{-1\} \\ S_3 = (-1, 1) \\ S_4 = \{1\} \\ S_5 = (1, \infty). \end{cases}$$

Above  $S_i$   $i = 1, 5$ , there are no semi-algebraic functions, and only one cell  $S_{i,1} = S_i \times \mathbb{R}$ .

Above  $S_i$ ,  $i = 2, 4$  there is only one semi-algebraic function associating to  $-1$  and  $1$  the constant value  $0$ , and there are three cells.

$$\begin{cases} S_{i,1} = S_i \times (-\infty, 0) \\ S_{i,2} = S_i \times \{0\} \\ S_{i,3} = S_i \times (0, \infty) \end{cases}$$

Above  $S_3$ , there are two semi-algebraic functions  $\xi_{3,1}$  and  $\xi_{3,2}$  associating to  $x \in S_3$  the values  $\xi_{3,1}(x) = -\sqrt{1-x^2}$  and  $\xi_{3,2}(x) = \sqrt{1-x^2}$ . There are 5 cells above  $S_3$ , the graphs of  $\xi_{3,1}$  and  $\xi_{3,2}$  and the bands they define

$$\begin{cases} S_{3,1} = \{(x, y) \mid -1 < x < 1, y < \xi_{3,1}(x)\} \\ S_{3,2} = \{(x, y) \mid -1 < x < 1, y = \xi_{3,1}(x)\} \\ S_{3,3} = \{(x, y) \mid -1 < x < 1, \xi_{3,1}(x) < y < \xi_{3,2}(x)\} \\ S_{3,4} = \{(x, y) \mid -1 < x < 1, y = \xi_{3,2}(x)\} \\ S_{3,5} = \{(x, y) \mid -1 < x < 1, \xi_{3,2}(x) < y\}. \end{cases}$$

Above  $S_{i,j}$ ,  $(i, j) \in \{(1, 1), (2, 1), (2, 3), (3, 1), (3, 5), (4, 1), (4, 3), (5, 1)\}$ , there are no semi-algebraic functions, and only one cell:

$$S_{i,j,1} = S_{i,j} \times \mathbb{R}$$

Above  $S_{i,j}$ ,  $(i, j) \in \{(2, 2), (3, 2), (3, 4), (4, 2)\}$ , there is only one semi-algebraic function, the constant function  $0$ , and three cells:

$$\begin{cases} S_{i,j,1} = S_{i,j} \times (-\infty, 0) \\ S_{i,j,2} = S_{i,j} \times \{0\} \\ S_{i,j,3} = S_{i,j} \times (0, \infty) \end{cases}$$

Above  $S_{3,3}$ , there are two semi-algebraic functions  $\xi_{3,3,1}$  and  $\xi_{3,3,2}$  associating to  $(x, y) \in S_{3,3}$  the values

$$\begin{aligned} \xi_{3,3,1}(x, y) &= -\sqrt{1-x^2-y^2} \\ \xi_{3,3,2}(x, y) &= \sqrt{1-x^2-y^2}, \end{aligned}$$

and five cells

$$\begin{cases} S_{3,3,1} = \{(x, y, z) \mid (x, y) \in S_{3,3}, z < \xi_{3,3,1}(x, y)\} \\ S_{3,3,2} = \{(x, y, z) \mid (x, y) \in S_{3,3}, z = \xi_{3,3,1}(x, y)\} \\ S_{3,3,3} = \{(x, y, z) \mid (x, y) \in S_{3,3}, \xi_{3,3,1}(x, y) < z < \xi_{3,3,2}(x, y)\} \\ S_{3,3,4} = \{(x, y, z) \mid (x, y) \in S_{3,3}, z = \xi_{3,3,2}(x, y)\} \\ S_{3,3,5} = \{(x, y, z) \mid (x, y) \in S_{3,3}, \xi_{3,3,2}(x, y) < z\}. \end{cases}$$

□

Note that a cylindrical decomposition has a recursive structure. Let  $S$  be a cell of level  $i$  of a cylindrical decomposition  $\mathcal{S}$  and  $x \in S$ . Denoting by  $\pi_i$  the canonical projection of  $\mathbb{R}^k$  to  $\mathbb{R}^i$  and identifying  $\pi_i^{-1}(x)$  with  $\mathbb{R}^{k-i}$ , the finite partition of  $\mathbb{R}^{k-i}$  obtained by intersecting the cells of  $\mathcal{S}_{i+j}$  above  $S$  with  $\pi_i^{-1}(x)$  is a cylindrical decomposition of  $\mathbb{R}^{k-i}$  called the **cylindrical decomposition induced by  $\mathcal{S}$  above  $x$** .

**Definition 5.5.** Given a finite set  $\mathcal{P}$  of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ , a subset  $S$  of  $\mathbb{R}^k$  is  **$\mathcal{P}$ -semi-algebraic** if  $S$  is the realization of a quantifier free formula with atoms  $P=0, P>0$  or  $P<0$  with  $P \in \mathcal{P}$ . It is clear that for every semi-algebraic subset  $S$  of  $\mathbb{R}^k$ , there exists a finite set  $\mathcal{P}$  of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$  such that  $S$  is  $\mathcal{P}$ -semi-algebraic.

A subset  $S$  of  $\mathbb{R}^k$  is  **$\mathcal{P}$ -invariant** if every polynomial  $P \in \mathcal{P}$  has a constant sign ( $>0, <0$ , or  $=0$ ) on  $S$ .

A **cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$**  is a cylindrical decomposition for which each cell  $C \in \mathcal{S}_k$  is  $\mathcal{P}$ -invariant. It is clear that if  $S$  is  $\mathcal{P}$ -semi-algebraic, a cylindrical decomposition adapted to  $\mathcal{P}$  is a cylindrical decomposition adapted to  $S$ . □

The main purpose of the next few pages is to prove the following result.

**Theorem 5.6. [Cylindrical decomposition]** *For every finite set  $\mathcal{P}$  of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ , there is a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .*

The theorem immediately implies:

**Corollary 5.7.** *For every finite family of semi-algebraic sets  $S_1, \dots, S_\ell$  of  $\mathbb{R}^k$ , there is a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $S_1, \dots, S_\ell$ .*

Since we intend to construct a cylindrical decomposition adapted to  $\mathcal{P}$  it is convenient if for  $S \in \mathcal{S}_{k-1}$  we choose each  $\xi_{S,j}$  to be a root of a polynomial  $P \in \mathcal{P}$ , as a function of  $(x_1, \dots, x_{k-1}) \in S$ . To this end, we shall prove that the real and complex roots (those in  $\mathbb{R}[i]=\mathbb{C}$ ) of a univariate polynomial depend continuously on its coefficients.

**Notation 5.8. [Disk]** We write  $D(z, r) = \{w \in \mathbb{C} \mid |z - w| < r\}$  for the **open disk** centered at  $z$  with radius  $r$ . □

First we need the following bound on the modulus of the roots of a polynomial.

**Proposition 5.9.** *Let  $P = a_p X^p + \dots + a_1 X + a_0 \in \mathbb{C}[X]$ , with  $a_p \neq 0$ . If  $x \in \mathbb{C}$  is a root of  $P$ , then*

$$|x| \leq \max_{i=1, \dots, p} \left( p \left| \frac{a_{p-i}}{a_p} \right| \right)^{1/i} = M .$$

**Proof:** If  $z \in \mathbb{C}$  is such that  $|z| > M$ , then  $|a_{p-i}| < |a_p| |z|^i / p, i = 1, \dots, p$ . Hence,

$$|a_{p-1} z^{p-1} + \dots + a_0| \leq |a_{p-1}| |z|^{p-1} + \dots + |a_0| < |a_p| |z|^p$$

and  $P(z) \neq 0$ . □

We identify the monic polynomial  $X^q + b_{q-1} X^{q-1} + \dots + b_0 \in \mathbb{C}[X]$  of degree  $q$  with the point  $(b_{q-1}, \dots, b_0) \in \mathbb{C}^q$ .



**Lemma 5.10.** *Given  $r > 0$ , there is an open neighborhood  $U$  of  $(X - z)^\mu$  in  $\mathbb{C}^\mu$  such that every monic polynomial in  $U$  has its roots in  $D(z, r)$ .*

**Proof:** Without loss of generality, we can suppose that  $z = 0$  and apply Proposition 5.9. □

Consider the map

$$\begin{aligned} m: \mathbb{C}^q \times \mathbb{C}^r &\longrightarrow \mathbb{C}^{q+r} \\ (Q, R) &\longmapsto QR \end{aligned}$$

defined by the multiplication of monic polynomials of degrees  $q$  and  $r$  respectively.

**Lemma 5.11.** *Let  $P_0 \in \mathbb{C}^{q+r}$  be a monic polynomial such that  $P_0 = Q_0 R_0$ , where  $Q_0$  and  $R_0$  are coprime monic polynomials of degrees  $q$  and  $r$ , respectively. There exist open neighborhoods  $U$  of  $P_0$  in  $\mathbb{C}^{q+r}$ ,  $U_1$  of  $Q_0$  in  $\mathbb{C}^q$  and  $U_2$  of  $R_0$  in  $\mathbb{C}^r$  such that any  $P \in U$  is uniquely the product of coprime monic polynomials  $QR$  with  $Q \in U_1$  and  $R \in U_2$ .*

**Proof:** The Jacobian matrix of  $m$  is the Sylvester matrix associated to

$$X^{q-1}R_0, \dots, R_0, X^{r-1}Q_0, \dots, Q_0$$

(Proposition 4.19). Hence the Jacobian of  $m$  is equal, up to sign, to the resultant of  $R_0$  and  $Q_0$  and is therefore different from zero by Proposition 4.15, since  $R_0$  and  $Q_0$  have no common factor. The conclusion follows using the implicit function theorem (Theorem 3.25). □

We can now prove

**Theorem 5.12. [Continuity of Roots]** *Let  $P \in \mathbb{R}[X_1, \dots, X_k]$  and let  $S$  be a semi-algebraic subset of  $\mathbb{R}^{k-1}$ . Assume that  $\deg(P(x', X_k))$  is constant on  $S$  and that for some  $a' \in S$ ,  $z_1, \dots, z_j$  are the distinct roots of  $P(a', X_k)$  in  $\mathbb{C} = \mathbb{R}[i]$ , with multiplicities  $\mu_1, \dots, \mu_j$ , respectively.*

*If the open disks  $D(z_i, r) \subset \mathbb{C}$  are disjoint then there is an open neighborhood  $V$  of  $a'$  such that for every  $x' \in V \cap S$ , the polynomial  $P(x', X_k)$  has exactly  $\mu_i$  roots, counted with multiplicity, in the disk  $D(z_i, r)$ , for  $i = 1, \dots, j$ .*

**Proof:** Let  $P_0 = (X - z_1)^{\mu_1} \dots (X - z_j)^{\mu_j}$ . By Lemma 5.11 there exist open neighborhoods  $U$  of  $P_0$  in  $\mathbb{C}^{\mu_1 + \dots + \mu_j}$ ,  $U_1$  of  $(X - z_1)^{\mu_1}$  in  $\mathbb{C}^{\mu_1}$ , ...,  $U_j$  of  $(X - z_j)^{\mu_j}$  in  $\mathbb{C}^{\mu_j}$  such that every  $P \in U$  can be factored uniquely as  $P = Q_1 \dots Q_j$ , where the  $Q_i$  are monic polynomials in  $U_i$ .

Using Lemma 5.10, it is clear that there is a neighborhood  $V$  of  $a'$  in  $S$  so that for every  $x' \in V$  the polynomial  $P(x', X_k)$  has exactly  $\mu_i$  roots counted with multiplicity in  $D(z_i, r)$ , for  $i = 1, \dots, j$ . □

We next consider the conditions which ensure that the zeros of two polynomials over a connected semi-algebraic set define a cylindrical structure.

**Proposition 5.13.** *Let  $P, Q \in \mathbb{R}[X_1, \dots, X_k]$  and  $S$  a semi-algebraically connected subset of  $\mathbb{R}^{k-1}$ . Suppose that  $P$  and  $Q$  are not identically 0 over  $S$  and that  $\deg(P(x', X_k)), \deg(Q(x', X_k)), \deg(\gcd(P(x', X_k), Q(x', X_k)))$ , the number of distinct roots of  $P(x', X_k)$  in  $\mathbb{C}$  and the number of distinct roots of  $Q(x', X_k)$  in  $\mathbb{C}$  are constant as  $x'$  varies over  $S$ . Then there exists  $\ell$  continuous semi-algebraic functions  $\xi_1 < \dots < \xi_\ell: S \rightarrow \mathbb{R}$  such that, for every  $x' \in S$ , the set of real roots of  $(PQ)(x', X_k)$  is exactly  $\{\xi_1(x'), \dots, \xi_\ell(x')\}$ .*

*Moreover for  $i = 1, \dots, \ell$ , the multiplicity of the root  $\xi_i(x')$  of  $P(x', X_k)$  (resp.  $Q(x', X_k)$ ) is constant for  $x' \in S$ . (If  $a$  is not a root, its multiplicity is zero, see Definition page 30).*

**Proof:** Let  $a' \in S$  and let  $z_1, \dots, z_j$  be the distinct roots in  $\mathbb{C}$  of the product  $(PQ)(a', X_k)$ . Let  $\mu_i$  (resp.  $\nu_i$ ) be the multiplicity of  $z_i$  as a root of  $P(a', X_k)$  (resp.  $Q(a', X_k)$ ). The degree of  $\gcd(P(a', X_k), Q(a', X_k))$  is  $\sum_{i=1}^j \min(\mu_i, \nu_i)$ , and each  $z_i$  has multiplicity  $\min(\mu_i, \nu_i)$  as a root of this greatest common divisor. Choose  $r > 0$  such that all disks  $D(z_i, r)$  are disjoint.

Using Theorem 5.12 and the fact that the number of distinct complex roots stays constant over  $S$ , we deduce that there exists a neighborhood  $V$  of  $a'$  in  $S$  such that for every  $x' \in V$ , each disk  $D(z_i, r)$  contains one root of multiplicity  $\mu_i$  of  $P(x', X_k)$  and one root of multiplicity  $\nu_i$  of  $Q(x', X_k)$ . Since the degree of  $\gcd(P(x', X_k), Q(x', X_k))$  is equal to  $\sum_{i=1}^j \min(\mu_i, \nu_i)$ , this greatest common divisor must have exactly one root  $\zeta_i$ , of multiplicity  $\min(\mu_i, \nu_i)$ , in each disk  $D(z_i, r)$  such that  $\min(\mu_i, \nu_i) > 0$ . So, for every  $x' \in V$ , and every  $i = 1, \dots, j$ , there is exactly one root  $\zeta_i$  of  $(PQ)(x', X_k)$  in  $D(z_i, r)$  which is a root of  $P(x', X_k)$  of multiplicity  $\mu_i$  and a root of  $Q(x', X_k)$  of multiplicity  $\nu_i$ . If  $z_i$  is real,  $\zeta_i$  is real (otherwise, its conjugate  $\bar{\zeta}_i$  would be another root of  $(PQ)(x', X_k)$  in  $D(z_i, r)$ ). If  $z_i$  is not real,  $\zeta_i$  is not real, since  $D(z_i, r)$  is disjoint from its image by conjugation. Hence, if  $x' \in V$ , the polynomial  $(PQ)(x', X_k)$  has the same number of distinct real roots as  $(PQ)(a', X_k)$ . Since  $S$  is semi-algebraically connected, the number of distinct real roots of  $(PQ)(x', X_k)$  is constant for  $x' \in S$  according to Proposition 3.9. Let  $\ell$  be this number. For  $1 \leq i \leq \ell$ , denote by  $\xi_i: S \rightarrow \mathbb{R}$  the function which sends  $x' \in S$  to the  $i$ -th real root (in increasing order) of  $(PQ)(x', X_k)$ . The argument above, with arbitrarily small  $r$  also shows that the functions  $\xi_i$  are continuous. It follows from the fact that  $S$  is semi-algebraically connected that each  $\xi_i(x')$  has constant multiplicity as a root of  $P(x', X_k)$  and as a root of  $Q(x', X_k)$  (cf Proposition 3.9). If  $S$  is described by the formula  $\Theta(X')$ , the graph of  $\xi_i$  is described by the formula

$$\begin{aligned} & \Theta(X') \\ & \wedge ( (\exists Y_1) \dots (\exists Y_\ell) (Y_1 < \dots < Y_\ell \wedge (PQ)(X', Y_1) = 0 \wedge \dots \wedge (PQ)(X', Y_\ell) = 0) \\ & \wedge ((\forall Y) (PQ)(X', Y) = 0 \Rightarrow (Y = Y_1 \vee \dots \vee Y = Y_\ell)) \wedge X_k = Y_i, \end{aligned}$$

which shows that  $\xi_i$  is semi-algebraic, by quantifier elimination (Corollary 2.78). □

We have thus proved:

**Proposition 5.14.** *Let  $\mathcal{P}$  be a finite subset of  $\mathbb{R}[X_1, \dots, X_k]$  and  $S$  a semi-algebraically connected semi-algebraic subset of  $\mathbb{R}^{k-1}$ . Suppose that, for every  $P \in \mathcal{P}$ ,  $\deg(P(x', X_k))$  and the number of distinct real roots of  $P$  are constant over  $S$  and that, for every pair  $P, Q \in \mathcal{P}$ ,  $\deg(\gcd(P(x', X_k), Q(x', X_k)))$  is also constant over  $S$ . Then there are  $\ell$  continuous semi-algebraic functions  $\xi_1 < \dots < \xi_\ell: S \rightarrow \mathbb{R}$  such that, for every  $x' \in S$ , the set of real roots of  $\prod_{P \in \mathcal{P}'} P(x', X_k)$ , where  $\mathcal{P}'$  is the subset of  $\mathcal{P}$  consisting of polynomials not identically 0 over  $S$ , is exactly  $\{\xi_1(x'), \dots, \xi_\ell(x')\}$ . Moreover, for  $i=1, \dots, \ell$  and for every  $P \in \mathcal{P}'$ , the multiplicity of the root  $\xi_i(x')$  of  $P(x', X_k)$  is constant for  $x' \in S$ .*

It follows from Chapter 4 (Proposition 4.24) that the number of distinct roots of  $P$ , of  $Q$  and the degree of the greatest common divisor of  $P$  and  $Q$  are determined by whether the signed subresultant coefficients  $\text{sRes}_i(P, P')$  and  $\text{sRes}_i(P, Q)$  are zero or not, as long as the degrees (with respect to  $X_k$ ) of  $P$  and  $Q$  are fixed.

**Notation 5.15. [Elimination]** Using Notation 1.16, with parameters  $X_1, \dots, X_{k-1}$  and main variable  $X_k$ , let

$$\text{Tru}(\mathcal{P}) = \{\text{Tru}(P) \mid P \in \mathcal{P}\}.$$

We define  $\text{Elim}_{X_k}(\mathcal{P})$  to be the set of polynomials in  $\mathbb{R}[X_1, \dots, X_{k-1}]$  defined as follows:

- If  $R \in \text{Tru}(\mathcal{P})$ ,  $\deg_{X_k}(R) \geq 2$ ,  $\text{Elim}_{X_k}(\mathcal{P})$  contains all  $\text{sRes}_j(R, \partial R / \partial X_k)$  which are not in  $\mathbb{R}$ ,  $j = 0, \dots, \deg_{X_k}(R) - 2$ .
- If  $R \in \text{Tru}(\mathcal{P})$ ,  $S \in \text{Tru}(\mathcal{P})$ ,
  - if  $\deg_{X_k}(R) > \deg_{X_k}(S)$ ,  $\text{Elim}_{X_k}(\mathcal{P})$  contains all  $\text{sRes}_j(R, S)$  which are not in  $\mathbb{R}$ ,  $j = 0, \dots, \deg_{X_k}(S) - 1$ ,
  - if  $\deg_{X_k}(R) < \deg_{X_k}(S)$ ,  $\text{Elim}_{X_k}(\mathcal{P})$  contains all  $\text{sRes}_j(S, R)$  which are not in  $\mathbb{R}$ ,  $j = 0, \dots, \deg_{X_k}(R) - 1$ ,
  - if  $\deg_{X_k}(R) = \deg_{X_k}(S)$ ,  $\text{Elim}_{X_k}(\mathcal{P})$  contains all  $\text{sRes}_j(S, \overline{R})$ , with  $\overline{R} = \text{lcof}(S)R - \text{lcof}(R)S$  which are not in  $\mathbb{R}$ ,  $j = 0, \dots, \deg_{X_k}(\overline{R}) - 1$ .
- If  $R \in \text{Tru}(\mathcal{P})$ , and  $\text{lcof}(R)$  is not in  $\mathbb{R}$ ,  $\text{Elim}_{X_k}(\mathcal{P})$  contains  $\text{lcof}(R)$ . □

**Theorem 5.16.** *Let  $\mathcal{P}$  be a set of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ , and let  $S$  be a semi-algebraically connected semi-algebraic subset of  $\mathbb{R}^{k-1}$  which is  $\text{Elim}_{X_k}(\mathcal{P})$ -invariant. Then there are continuous semi-algebraic functions  $\xi_1 < \dots < \xi_\ell: S \rightarrow \mathbb{R}$  such that, for every  $x' \in S$ , the set  $\{\xi_1(x'), \dots, \xi_\ell(x')\}$  is the set of all real roots of all non-zero polynomials  $P(x', X_k)$ ,  $P \in \mathcal{P}$ . The graph of each  $\xi_i$  (resp. each band of the cylinder  $S \times \mathbb{R}$  bounded by these graphs) is a semi-algebraically connected semi-algebraic set semi-algebraically homeomorphic to  $S$  (resp.  $S \times (0, 1)$ ) and is  $\mathcal{P}$ -invariant.*

**Proof:** For  $P$  in  $\mathcal{P}$ ,  $R \in \text{Tru}(P)$ , consider the constructible set  $A \subset \mathbb{R}^{k-1}$  defined by  $\text{lcof}(R) \neq 0, \deg(P) = \deg(R)$ . By Proposition 4.24, for every  $a' \in A$ , the vanishing or non-vanishing of the  $\text{sRes}_j(R, \partial R/\partial X_k)(a')$  determines the number of distinct roots of  $P(a', X_k)$  in  $\mathbb{C}$ , which is

$$\deg(R(a', X_k)) - \deg(\gcd(R(a', X_k), \partial R/\partial X_k(a', X_k)))$$

Similarly, for  $R \in \text{Tru}(P), S \in \text{Tru}(Q)$ , consider the constructible set  $B$  defined by

$$\text{lcof}(R) \neq 0, \deg(P) = \deg(R), \text{lcof}(S) \neq 0, \deg(Q) = \deg(S).$$

For every  $a' \in B$ , which of the  $\text{sRes}_j(R, S)(a')$  (resp.  $\text{sRes}_j(S, R)(a')$ ,  $\text{sRes}_j(S, \bar{R})(a')$ ) vanish, determine  $\deg(\gcd(P(a', X_k), Q(a', X_k)))$ , by Proposition 4.24. Thus, the assumption that a connected semi-algebraic subset of  $\mathbb{R}^{k-1}$  is  $\text{Elim}_{X_k}(\mathcal{P})$ -invariant implies that the hypotheses of Proposition 5.14 are satisfied. □

We are finally ready for the proof of Theorem 5.6.

**Proof of Theorem 5.6** The proof is by induction on the dimension of the ambient space.

Let  $\mathcal{Q} \subset \mathbb{R}[X_1]$  be finite. It is clear that there is a cylindrical decomposition of  $\mathbb{R}$  adapted to  $\mathcal{Q}$  since the real roots of the polynomials in  $\mathcal{Q}$  decompose the line into finitely many points and open intervals which constitute the cells of a cylindrical decomposition of  $\mathbb{R}$  adapted to  $\mathcal{Q}$ .

Let  $\mathcal{Q} \subset \mathbb{R}[X_1, \dots, X_i]$  be finite. Starting from a cylindrical decomposition of  $\mathbb{R}^{i-1}$  adapted to  $\text{Elim}_{X_i}(\mathcal{Q})$ , and applying to the cells of this cylindrical decomposition Proposition 5.16, yields a cylindrical decomposition of  $\mathbb{R}^i$  adapted to  $\mathcal{Q}$ .

This proves the theorem. □

*Example 5.17.* We illustrate this result by presenting a cylindrical decomposition of  $\mathbb{R}^3$  adapted to the polynomial  $P = X_1^2 + X_2^2 + X_3^2 - 1$ . The 0-th Sylvester-Habicht matrix of  $P$  and  $\partial P/\partial X_3$  is

$$\begin{bmatrix} 1 & 0 & X_1^2 + X_2^2 - 1 \\ 0 & 2 & 0 \\ 2 & 0 & 0 \end{bmatrix}.$$

Hence,  $\text{sRes}_0(P, \partial P/\partial X_3) = -4(X_1^2 + X_2^2 - 1)$  and  $\text{sRes}_1(P, \partial P/\partial X_3) = 2$ . Getting rid of irrelevant constant factors, we obtain

$$\text{Elim}_{X_3}(P) = \{X_1^2 + X_2^2 - 1\}.$$

Similarly,

$$\text{Elim}_{X_2}(\text{Elim}_{X_3}(P)) = \{X_1^2 - 1\}.$$

The associated cylindrical decomposition is precisely the one described in Example 5.4. □

*Remark 5.18.* The proof of Theorem 5.6 provides a method for constructing a cylindrical decomposition adapted to  $\mathcal{P}$ . In a projection phase, we eliminate the variables one after the other, by computing  $\text{Elim}_{X_k}(\mathcal{P})$ , then  $\text{Elim}_{X_{k-1}}(\text{Elim}_{X_k}(\mathcal{P}))$  etc. until we obtain a finite family of univariate polynomials.

In a lifting phase, we decompose the line in a finite number of cells which are the points and intervals defined by the family of univariate polynomials. Then we decompose the cylinder contained in  $\mathbb{R}^2$  above each of these points and intervals in a finite number of cells consisting of graphs and bands between these graphs. Then we decompose the cylinder contained in  $\mathbb{R}^2$  above each of plane cells in a finite number of cells consisting of graphs and bands between these graphs etc.

Note that the projection phase of the construction provides in fact an algorithm computing explicitly a family of polynomials in one variable. The complexity of this algorithm will be studied in Chapter 12.  $\square$

**Theorem 5.19.** *Every semi-algebraic subset  $S$  of  $\mathbb{R}^k$  is the disjoint union of a finite number of semi-algebraic sets, each of them semi-algebraically homeomorphic to an open  $i$ -cube  $(0, 1)^i \subset \mathbb{R}^i$  for some  $i \leq k$  (by convention  $(0, 1)^0$  is a point).*

**Proof:** According to Corollary 5.7, there exists a cylindrical decomposition adapted to  $S$ . Since these cells are homeomorphic to an open  $i$ -cube  $(0, 1)^i \subset \mathbb{R}^i$  for some  $i \leq k$ , the conclusion follows immediately.  $\square$

An easy consequence is the following which asserts the piecewise continuity of semi-algebraic functions.

**Proposition 5.20.** *Let  $S$  be a semi-algebraic set and let  $f: S \rightarrow \mathbb{R}^k$  be a semi-algebraic function. There is a partition of  $S$  in a finite number of semi-algebraic sets  $S_1, \dots, S_n$  such that the restriction  $f_i$  of  $f$  to  $S_i$  is semi-algebraic and continuous.*

**Proof:** By Theorem 5.19, the graph  $G$  of  $f$  is the union of open  $i$ -cubes of various dimensions, which are clearly the graphs of semi-algebraic continuous functions.  $\square$

## 5.2 Semi-algebraically Connected Components

**Theorem 5.21.** *Every semi-algebraic set  $S$  of  $\mathbb{R}^k$  is the disjoint union of a finite number of semi-algebraically connected semi-algebraic sets  $C_1, \dots, C_\ell$  that are both closed and open in  $S$ .*

The  $C_1, \dots, C_\ell$  are called the **semi-algebraically connected components** of  $S$ .

**Proof of Theorem 5.21:** By Theorem 5.19,  $S$  is the disjoint union of a finite number of semi-algebraic sets  $S_i$  semi-algebraically homeomorphic to open  $d(i)$ -cubes  $(0, 1)^{d(i)}$  and hence semi-algebraically connected by Proposition 3.8. Consider the smallest equivalence relation  $\mathcal{R}$  on the set of the  $S_i$  containing the relation “ $S_i \cap \overline{S_j} \neq \emptyset$ ”. Let  $C_1, \dots, C_\ell$  be the unions of the equivalence classes for  $\mathcal{R}$ . The  $C_j$  are semi-algebraic, disjoint, closed in  $S$ , and their union is  $S$ . We show now that each  $C_j$  is semi-algebraically connected. Suppose that we have  $C_j = F_1 \cup F_2$  with  $F_1$  and  $F_2$  disjoint, semi-algebraic and closed in  $C_j$ . Since each  $S_i$  is semi-algebraically connected,  $S_i \subset C_j$  implies that  $S_i \subset F_1$  or  $S_i \subset F_2$ . Since  $F_1$  (resp.  $F_2$ ) is closed in  $C_j$ , if  $S_i \subset F_1$  (resp.  $F_2$ ) and  $S_i \cap \overline{S_{i'}} \neq \emptyset$  then  $S_{i'} \subset F_1$  (resp.  $F_2$ ). By the definition of the  $C_j$ , we have  $C_j = F_1$  or  $C_j = F_2$ . So  $C_j$  is semi-algebraically connected.  $\square$

**Theorem 5.22.** *A semi-algebraic subset  $S$  of  $\mathbb{R}^k$  is semi-algebraically connected if and only if it is connected. Every semi-algebraic set (and in particular every algebraic subset) of  $\mathbb{R}^k$  has a finite number of connected components, each of which is semi-algebraic.*

**Proof:** It is clear that if  $S$  is connected, it is semi-algebraically connected.

If  $S$  is not connected then there exist open sets  $O_1$  and  $O_2$  (not necessarily semi-algebraic) with

$$S \subset O_1 \cup O_2, O_1 \cap S \neq \emptyset, O_2 \cap S \neq \emptyset$$

and  $(S \cap O_1) \cap (S \cap O_2) = \emptyset$ . By Theorem 5.19, we know that  $S$  is a union of a finite number  $C_1, \dots, C_\ell$  of semi-algebraic sets homeomorphic to open cubes of various dimensions. If  $O_1 \cap S$  and  $O_2 \cap S$  are unions of a finite number of semi-algebraic sets among  $C_1, \dots, C_\ell$ ,  $O_1 \cap S$  and  $O_2 \cap S$  are semi-algebraic and  $S$  is not semi-algebraically connected. Otherwise, some  $C_i$  is disconnected by  $O_1$  and  $O_2$ , which is impossible since  $C_i$  is homeomorphic to an open cube.

Hence a semi-algebraic subset  $S$  of  $\mathbb{R}^k$  is semi-algebraically connected if and only if it is connected. The remainder of the theorem follows from Theorem 5.21.  $\square$

**Theorem 5.23.** *A semi-algebraic set is semi-algebraically connected if and only if it is semi-algebraically path connected.*

**Proof:** Since  $[0, 1]$  is semi-algebraically connected, it is clear that semi-algebraic path connectedness implies semi-algebraic connectedness. We prove the converse by using Theorem 5.19 and the proof of Theorem 5.21. It is obvious that an open  $d$ -cube is semi-algebraically path connected. It is then enough to show that if  $S_i$  and  $S_j$  are semi-algebraically homeomorphic to open  $d$ -cubes, with  $S_i \cap \overline{S_j} \neq \emptyset$ , then  $S_i \cup S_j$  is semi-algebraically path connected. But this is a straightforward consequence of the Curve Selection Lemma (Theorem 3.19).  $\square$

Let  $R'$  be a real closed extension of the real closed field  $R$ .

**Proposition 5.24.** *The semi-algebraic set  $S$  is semi-algebraically connected if and only if  $\text{Ext}(S, \mathbb{R}')$  is semi-algebraically connected.*

*More generally, if  $C_1, \dots, C_\ell$  are the semi-algebraically connected components of  $S$ , then  $\text{Ext}(C_1, \mathbb{R}'), \dots, \text{Ext}(C_\ell, \mathbb{R}')$  are the semi-algebraically connected components of  $\text{Ext}(S, \mathbb{R}')$ .*

**Proof:** Given a decomposition  $S = \bigcup_{i=1}^m S_i$ , with, for each  $i$ , a semi-algebraic homeomorphism  $\varphi_i: (0, 1)^{d(i)} \rightarrow S_i$ , the extension gives a decomposition

$$\text{Ext}(S, \mathbb{R}') = \bigcup_{i=1}^m \text{Ext}(S_i, \mathbb{R}'),$$

and semi-algebraic homeomorphisms

$$\text{Ext}(\varphi_i, \mathbb{R}'): (\text{Ext}((0, 1), \mathbb{R}')^{d_i} \rightarrow (\text{Ext}(S_i, \mathbb{R}')).$$

The characterization of the semi-algebraically connected components from a decomposition (cf. Theorem 5.21) then gives the result.  $\square$

### 5.3 Dimension

Let  $S$  be a semi-algebraic subset of  $\mathbb{R}^k$ . Take a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $S$ . A naive definition of the dimension of  $S$  is the maximum of the dimension of the cells contained in  $S$ , the dimension of a cell semi-algebraically homeomorphic to  $(0, 1)^d$  being  $d$ . But this definition is not intrinsic. We would have to prove that the dimension so defined does not depend on the choice of a cylindrical decomposition adapted to  $S$ . Instead, we introduce an intrinsic definition of dimension and show that it coincides with the naive one.

The **dimension**  $\dim(S)$  of a semi-algebraic set  $S$  is the largest  $d$  such that there exists an injective semi-algebraic map from  $(0, 1)^d$  to  $S$ . By convention, the dimension of the empty set is  $-1$ . Note that the dimension of a set is clearly invariant under semi-algebraic bijections. Observe that it is not obvious for the moment that the dimension is always  $< +\infty$ . It is also not clear that this definition of dimension agrees with the intuitive notion of dimension for cells.

We are going to prove the following result.

**Theorem 5.25.** *Let  $S \subset \mathbb{R}^k$  be semi-algebraic and consider a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $S$ . Then the dimension of  $S$  is finite and is the maximum dimension of the cells contained in  $S$ .*

The key ingredient for proving this result is the following lemma.

**Lemma 5.26.** *Let  $S$  be a semi-algebraic subset of  $\mathbb{R}^k$  with non-empty interior. Let  $f: S \rightarrow \mathbb{R}^k$  be an injective semi-algebraic map. Then  $f(S)$  has non-empty interior.*

**Proof:** We prove the lemma by induction on  $k$ . If  $k = 1$ ,  $S$  is semi-algebraic and has infinite cardinality, hence  $f(S) \subset \mathbb{R}$  is semi-algebraic and infinite and must therefore contain an interval.

Assume that  $k > 1$  and that the lemma is proved for all  $\ell < k$ . Using the piecewise continuity of semi-algebraic functions (Proposition 5.20), we can assume moreover that  $f$  is continuous. Take a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $f(S)$ . If  $f(S)$  has empty interior, it contains no cell open in  $\mathbb{R}^k$ . Hence  $f(S)$  is the union of cells  $C_1, \dots, C_n$  that are not open in  $\mathbb{R}^k$  and, for  $i = 1, \dots, n$ , there is a semi-algebraic homeomorphism  $C_i \rightarrow (0, 1)^{\ell_i}$  with  $\ell_i < k$ . Take a cylindrical decomposition of  $\mathbb{R}^k$  adapted to the  $f^{-1}(C_i)$ . Since  $S = \bigcup_{i=1}^n f^{-1}(C_i)$  has non-empty interior, one of the  $f^{-1}(C_i)$ , say  $f^{-1}(C_1)$ , must contain a cell  $C$  open in  $\mathbb{R}^k$ . The restriction of  $f$  to  $C$  gives an injective continuous semi-algebraic map  $C \rightarrow C_1$ .

Since  $C$  is semi-algebraically homeomorphic to  $(0, 1)^k$  and  $C_1$  semi-algebraically homeomorphic to  $(0, 1)^\ell$  with  $\ell < k$ , we obtain an injective continuous semi-algebraic map  $g$  from  $(0, 1)^k$  to  $(0, 1)^\ell$ . Set  $a = (\frac{1}{2}, \dots, \frac{1}{2}) \in \mathbb{R}^{k-\ell}$  and consider the mapping  $g_a$  from  $(0, 1)^\ell$  to  $(0, 1)^\ell$  defined by  $g_a(x) = g(a, x)$ . We can apply the inductive assumption to  $g_a$ . It implies that  $g_a((0, 1)^\ell)$  has non-empty interior in  $(0, 1)^\ell$ . Choose a point  $c = g_a(b)$  in the interior of  $g_a((0, 1)^\ell)$ . Since  $g$  is continuous, all points close enough to  $(a, b)$  are mapped by  $g$  to the interior of  $g_a((0, 1)^\ell)$ . Let  $(x, b)$  be such a point with  $x \neq a$ . Since  $g_a$  is onto the interior of  $g_a((0, 1)^\ell)$  there is  $y \in (0, 1)^\ell$  such that  $g(x, b) = g_a(y) = g(a, y)$ , which contradicts the fact that  $g$  is injective. Hence,  $f(S)$  has non-empty interior.  $\square$

**Proposition 5.27.** *The dimension of  $(0, 1)^d$  is  $d$ . The dimension of a cell semi-algebraically homeomorphic to  $(0, 1)^d$  is  $d$ .*

**Proof:** There is no injective semi-algebraic map from  $(0, 1)^e$  to  $(0, 1)^d$  if  $e > d$ . Otherwise, the composition of such a map with the embedding of  $(0, 1)^d$  in  $\mathbb{R}^e = \mathbb{R}^d \times \mathbb{R}^{e-d}$  as  $(0, 1)^d \times \{0\}$  would contradict Lemma 5.26. This shows the first part of the corollary. The second part follows, using the fact that the dimension, according to its definition, is invariant under semi-algebraic bijection.  $\square$

**Proposition 5.28.** *If  $S \subset T$  are semi-algebraic sets,  $\dim(S) \leq \dim(T)$ .*

*If  $S$  and  $T$  are semi-algebraic subsets of  $\mathbb{R}^k$ ,  $\dim(S \cup T) = \max(\dim(S), \dim(T))$ .*

*If  $S$  and  $T$  are semi-algebraic sets,  $\dim(S \times T) = \dim(S) + \dim(T)$ .*

**Proof:** That  $\dim(S) \leq \dim(T)$  is clear from the definition. The inequality  $\dim(S \cup T) \geq \max(\dim S, \dim T)$  follows from 1. Now let  $f: (0, 1)^d \rightarrow S \cup T$  be a semi-algebraic injective map. Taking a cylindrical decomposition of  $\mathbb{R}^d$  adapted to  $f^{-1}(S)$  and  $f^{-1}(T)$ , we see that  $f^{-1}(S)$  or  $f^{-1}(T)$  contains a cell of dimension  $d$ . Since  $f$  is injective, we have  $\dim(S) \geq d$  or  $\dim(T) \geq d$ . This proves the reverse inequality  $\dim(S \cup T) \leq \max(\dim(S), \dim(T))$ .



Since  $\dim(S \cup T) = \max(\dim(S), \dim(T))$ , it is sufficient to consider the case where  $S$  and  $T$  are cells.

Since  $S \times T$  is semi-algebraically homeomorphic to  $(0, 1)^{\dim S} \times (0, 1)^{\dim T}$ , the assertion in this case follows from Proposition 5.27.  $\square$

**Proof of Theorem 5.25:** The result follows immediately from Proposition 5.27 and Proposition 5.28.  $\square$

**Proposition 5.29.** *Let  $S$  be a semi-algebraic subset of  $\mathbb{R}^k$ , and let  $f: S \rightarrow \mathbb{R}^\ell$  a semi-algebraic mapping. Then  $\dim(f(S)) \leq \dim(S)$ . If  $f$  is injective, then  $\dim(f(S)) = \dim(S)$ .*

The proof uses the following lemma.

**Lemma 5.30.** *Let  $S \subset \mathbb{R}^{k+\ell}$  be a semi-algebraic set,  $\pi$  the projection of  $\mathbb{R}^{k+\ell}$  onto  $\mathbb{R}^\ell$ . Then  $\dim(\pi(S)) \leq \dim(S)$ . If, moreover, the restriction of  $\pi$  to  $S$  is injective, then  $\dim(\pi(S)) = \dim S$ .*

**Proof:** When  $\ell = 1$  and  $S$  is a graph or a band in a cylindrical decomposition of  $\mathbb{R}^{k+1}$ , the result is clear. If  $S$  is any semi-algebraic subset of  $\mathbb{R}^{k+1}$ , it is a union of such cells for a decomposition, and the result is still true. The case of any  $\ell$  follows by induction.  $\square$

**Proof of Proposition 5.30:** Let  $G \subset \mathbb{R}^{k+\ell}$  be the graph of  $f$ . Lemma 5.30 tells us that  $\dim(S) = \dim(G)$  and  $\dim(f(S)) \leq \dim(S)$ , with equality if  $f$  is injective.  $\square$

Finally the following is clear:

**Proposition 5.31.** *Let  $V$  be an  $S^\infty$  submanifold of dimension  $d$  of  $\mathbb{R}^k$  (as a submanifold of  $\mathbb{R}^k$ , see Definition 3.25). Then the dimension of  $V$  as a semi-algebraic set is  $d$ .*

## 5.4 Semi-algebraic Description of Cells

In the preceding sections, we decomposed semi-algebraic sets into simple pieces, the cells, which are semi-algebraically homeomorphic to open  $i$ -cubes. We have also explained how to produce such a decomposition adapted to a finite set of polynomials  $\mathcal{P}$ . But the result obtained is not quite satisfactory, as we do not have a semi-algebraic description of the cells by a boolean combination of polynomial equations and inequalities. Since the cells are semi-algebraic, this description certainly exists. It would be nice to have the polynomials defining the cells of a cylindrical decomposition adapted to  $\mathcal{P}$ . This will be possible with the help of the derivatives of the polynomials.

We need to introduce a few definitions.

**Definition 5.32.** A **weak sign condition** is an element of

$$\{\{0\}, \{0, 1\}, \{0, -1\}\}.$$

Note that

$$\begin{cases} \text{sign}(x) \in \{0\} & \text{if and only if } x = 0, \\ \text{sign}(x) \in \{0, 1\} & \text{if and only if } x \geq 0, \\ \text{sign}(x) \in \{0, -1\} & \text{if and only if } x \leq 0. \end{cases}$$

A **weak sign condition** on  $\mathcal{Q}$  is an element of  $\{\{0\}, \{0, 1\}, \{0, -1\}\}^{\mathcal{Q}}$ . If  $\sigma \in \{0, 1, -1\}^{\mathcal{Q}}$ , its **relaxation**  $\bar{\sigma}$  is the weak sign condition on  $\mathcal{Q}$  defined by  $\bar{\sigma}(Q) = \overline{\sigma(Q)}$ . The **realization of the weak sign condition**  $\tau$  is

$$\text{Reali}(\tau) = \{x \in \mathbb{R}^k \mid \bigwedge_{Q \in \mathcal{Q}} \text{sign}(Q(x)) \in \tau(Q)\}.$$

The weak sign condition  $\tau$  is **realizable** if  $\text{Reali}(\tau)$  is non-empty.

A set of polynomials  $\mathcal{Q} \subset \mathbb{R}[X]$  is **closed under differentiation** if  $0 \notin \mathcal{Q}$  and if for each  $Q \in \mathcal{Q}$  then  $Q' \in \mathcal{Q}$  or  $Q' = 0$ .  $\square$

The following result is an extension of the Basic Thom's lemma (Lemma 2.28) seen in Chapter 2. It implies that if a family of polynomials is stable under differentiation, the cells it defines on a line are described by sign conditions on this family.

**Lemma 5.33. [Thom's lemma]** *Let  $\mathcal{Q} \subset \mathbb{R}[X]$  be a finite set of polynomials closed under differentiation and let  $\sigma$  be a sign condition on the set  $\mathcal{Q}$ . Then*

- $\text{Reali}(\sigma)$  is either empty, a point, or an open interval.
- If  $\text{Reali}(\sigma)$  is empty, then  $\text{Reali}(\bar{\sigma})$  is either empty or a point.
- If  $\text{Reali}(\sigma)$  is a point, then  $\text{Reali}(\bar{\sigma})$  is the same point.
- If  $\text{Reali}(\sigma)$  is an open interval then  $\text{Reali}(\bar{\sigma})$  is the corresponding closed interval.

**Proof:** The proof is by induction on  $s$ , the number of polynomials in  $\mathcal{Q}$ . There is nothing to prove if  $s = 0$ . Suppose that the proposition has been proved for  $s$  and that  $Q$  has maximal degree in  $\mathcal{Q}$ , which is closed under differentiation and has  $s + 1$  elements. The set  $\mathcal{Q} \setminus \{Q\}$  is also closed under differentiation. Let  $\sigma \in \{0, 1, -1\}^{\mathcal{Q}}$  be a sign condition on  $\mathcal{Q}$ , and let  $\sigma'$  be its restriction to  $\mathcal{Q} \setminus \{Q\}$ . If  $\text{Reali}(\sigma')$  is either a point or empty, then

$$\text{Reali}(\sigma) = \text{Reali}(\sigma') \cap \{x \in \mathbb{R} \mid \text{sign}(Q(x)) = \sigma(Q)\}$$

is either a point or empty. If  $\text{Reali}(\sigma')$  is an open interval, the derivative of  $Q$  (which is among  $\mathcal{Q} \setminus \{Q\}$ ), has a constant non-zero sign on it (except if  $Q$  is a constant, which is a trivial case). Thus  $Q$  is strictly monotone on  $\text{Reali}(\bar{\sigma}')$  so that the claimed properties are satisfied for  $\text{Reali}(\sigma)$ .  $\square$

By alternately applying the operation Elim and closing under differentiation we obtain a set of polynomials whose realizable sign conditions define the cells of a cylindrical decomposition adapted to  $\mathcal{P}$ .

**Theorem 5.34.** *Let  $\mathcal{P}^* = \cup_{i=1, \dots, k} \mathcal{P}_i$  be a finite set of non-zero polynomials such that:*

- $\mathcal{P}_k$  contains  $\mathcal{P}$ ,
- for each  $i$ ,  $\mathcal{P}_i$  is a subset of  $\mathbb{R}[X_1, \dots, X_i]$  that is closed under differentiation with respect to  $X_i$ ,
- for  $i \leq k$ ,  $\text{Elim}_{X_i}(\mathcal{P}_i) \subset \mathcal{P}_{i-1}$ .

Writing  $\mathcal{P}_{\leq i} = \bigcup_{j \leq i} \mathcal{P}_j$ , the families  $\mathcal{S}_i$ , for  $i = 1, \dots, k$ , consisting of all  $\text{Reali}(\sigma)$  with  $\sigma$  a realizable sign condition on  $\mathcal{P}_{\leq i}$  constitute a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .

**Proof:** The case  $k = 1$  is covered by Lemma 5.33. The proof of the general case is clear by induction on  $k$ , again using Lemma 5.33.  $\square$

*Remark 5.35.* Since  $\mathcal{P}_{\leq i+1}$  is closed under differentiation, for every cell  $S \subset \mathbb{R}^i$  and every semi-algebraic function  $\xi_{S,j}$  of the cylindrical decomposition described in the theorem, there exists  $P \in \mathcal{P}_{\leq i+1}$  such that, for every  $x \in S$ ,  $\xi_{S,j}(x)$  is a simple root of  $P(x, X_{i+1})$ .  $\square$

## 5.5 Stratification

We do not have so far much information concerning which cells of a cylindrical decomposition are adjacent to others, for cells which are not above the same cell.

In the case of the cylindrical decomposition adapted to the sphere, it is not difficult to determine the topology of the sphere from the cell decomposition. Indeed, the two functions on the disk defined by  $X_1^2 + X_2^2 + X_3^2 < 1$ , whose graphs are the two open hemispheres, have an obvious extension by continuity to the closed disk.

*Example 5.36.* We give an example of a cylindrical decomposition where it is not the case that the functions defined on the cells have an extension by continuity to boundary of the cell. Take  $P = (X_1 X_2 X_3) - (X_1^2 + X_2^2)$ . In order to visualize the corresponding zero set, it is convenient to fix the value of  $x_3$ .

The zero set  $Z$  of  $(X_1 X_2 x_3) - (X_1^2 + X_2^2)$  can be described as follows.

- If  $-2 < x_3 < 2$ ,  $Z$  consists of the isolated point  $(0, 0, x_3)$ .
- If  $x_3 = -2$  or  $x_3 = 2$ ,  $Z$  consists of one double line through the origin in the plane  $X_3 = x_3$ .

- If  $x_3 > -2$  or  $x_3 < 2$ ,  $Z$  consists of two distinct lines through the origin in the plane  $X_3 = x_3$ .

Note that the set of zeroes of  $P$  in the ball of center 0 and radius of 1 is the segment  $(-1, 1)$  of the  $X_3$  axis, so that  $\text{Zer}(P, \mathbb{R}^3)$  has an open subset which is a semi-algebraic set of dimension 1.

- When  $X_1 X_2 \neq 0$ ,  $P = 0$  is equivalent to

$$X_3 = \frac{X_1^2 + X_2^2}{X_1 X_2}.$$

- When  $X_1 = 0, X_2 \neq 0$ , the polynomial  $P$  is  $-X_2^2$ .
- When  $X_2 = 0, X_1 \neq 0$ , the polynomial  $P$  is  $-X_1^2$ .
- When  $X_1 = 0, X_2 = 0$ ,  $P$  is identically zero.

The function  $(X_1^2 + X_2^2)/(X_1 X_2)$  does not have a well defined limit when  $X_1$  and  $X_2$  tend to 0. The function describing the zeros of  $P$  on each open quadrant cannot be extended continuously to the closed quadrant.

The main difference with the example of the sphere is the fact that the polynomial  $P$  is not monic as polynomial in  $X_3$ : the leading coefficient  $X_1 X_2$  vanishes, and  $P$  is even identically zero for  $X_1 = X_2 = 0$ . □

We explain now that the information provided by the cylindrical decomposition is not sufficient to determine the topology.

*Example 5.37.* We describe two surfaces having the same cylindrical decomposition and a different topology, namely the two surfaces defined as the zero sets of

$$\begin{aligned} P_1 &= (X_1 X_2 X_3)^2 - (X_1^2 + X_2^2)^2 \\ P_2 &= P_2 = (X_1 X_2 X_3 - (X_1 - X_2)^2)(X_1 X_2 X_3 - (X_1 + X_2)^2). \end{aligned}$$

Consider first  $P_1 = (X_1 X_2 X_3)^2 - (X_1^2 + X_2^2)^2$ . In order to visualize the zero set of  $P_1$ , it is convenient to fix the value of  $x_3$ .

The zero set of  $P_1 = (X_1 X_2 x_3)^2 - (X_1^2 + X_2^2)^2$  is the union of the zero set  $Z_1$  of  $(X_1 X_2 x_3) + (X_1^2 + X_2^2)$  and the zero set  $Z_2$  of  $(X_1 X_2 x_3) - (X_1^2 + X_2^2)$  in the plane  $X_3 = x_3$ .

- If  $-2 < x_3 < 2$ ,  $Z_1 = Z_2$  consists of the isolated point  $(0, 0, x_3)$ .
- If  $x_3 = -2$  or  $x_3 = 2$ ,  $Z_1 \cup Z_2$  consists of two distinct lines through the origin in the plane  $X_3 = x_3$ .
- If  $x_3 > -2$  or  $x_3 < 2$ ,  $Z_1 \cup Z_2$  consists of four distinct lines through the origin in the plane  $X_3 = x_3$ .

Note that the set of zeroes of  $P_1$  in the ball of center 0 and radius of 1 is the segment  $(-1, 1)$  of the  $X_3$  axis, so that  $\text{Zer}(P_1, \mathbb{R}^3)$  has an open subset which is a semi-algebraic set of dimension 1.

It is easy to see that the 9 cells of  $\mathbb{R}^2$  defined by the signs of  $X_1$  and  $X_2$  together with the 3 cells of  $\mathbb{R}$  defined by the sign of  $X_1$  determine a cylindrical decomposition adapted to  $\{P_1\}$ .

- When  $X_1 X_2 \neq 0$ ,  $P_1 = 0$  is equivalent to

$$X_3 = \frac{X_1^2 + X_2^2}{X_1 X_2} \quad \text{or} \quad X_3 = -\frac{X_1^2 + X_2^2}{X_1 X_2}.$$

So the zeroes of  $P_1$  are described by two graphs of functions over each open quadrant, and the cylindrical decomposition of  $P_1$  has five cells over each open quadrant. The sign of  $P_1$  in these five cells is  $1, 0, -1, 0, 1$ .

- When  $X_1 = 0$ ,  $X_2 \neq 0$ , the polynomial  $P_1$  is  $-X_2^4$ . The cylinders over each open half-axis have one cell on which  $P_1$  is negative.
- When  $X_1 \neq 0$ ,  $X_2 = 0$ , the polynomial  $P_1$  is  $-X_1^4$ . The cylinders over each open half-axis have one cell on which  $P_1$  is negative.
- When  $X_1 = 0$ ,  $X_2 = 0$ ,  $P_1$  is identically zero. The cylinder over the origin has one cell on which  $P_1$  is zero.

The function  $(X_1^2 + X_2^2)/(X_1 X_2)$  does not have a well defined limit when  $X_1$  and  $X_2$  tend to 0. Moreover, the closure of the graph of the function  $(X_1^2 + X_2^2)/(X_1 X_2)$  on  $X_1 > 0$ ,  $X_2 > 0$  intersected with the line above the origin is  $[2, +\infty)$ , which is not a cell of the cylindrical decomposition.

Consider now  $P_2 = (X_1 X_2 X_3 - (X_1 - X_2)^2)(X_1 X_2 X_3 - (X_1 + X_2)^2)$ ,

In order to visualize the corresponding zero set, it is convenient to fix the value of  $x_3$ .

The zero set of  $(X_1 X_2 x_3 - (X_1 - X_2)^2)(X_1 X_2 x_3 - (X_1 + X_2)^2)$  is the union of the zero set  $Z_1$  of  $X_1 X_2 x_3 - (X_1 - X_2)^2$  and the zero set  $Z_2$  of  $X_1 X_2 x_3 - (X_1 + X_2)^2$  in the plane  $X_3 = x_3$ .

It can be easily checked that:

- If  $-4 < x_3$ , or  $x_3 > 4$ , the zeroes of  $P_2$  in the plane  $X_3 = x_3$  consist of four lines through the origin.
- If  $x_3 = -4$  or  $x_3 = 4$ , the zeroes of  $P_2$  in the plane  $X_3 = x_3$  consists of three lines through the origin.
- If  $x_3 = 0$ , the zeroes of  $P_2$  in the plane  $X_3 = x_3$  consists of two lines through the origin.
- If  $-4 < x_3 < 0$  or  $0 < x_3 < 4$ , the zeroes of  $P_2$  in the plane  $X_3 = x_3$  consists of two lines through the origin.

It is also easy to see that the 9 cells of  $\mathbb{R}^2$  defined by the signs of  $X_1$  and  $X_2$  and the 3 cells of  $\mathbb{R}$  defined by the sign of  $X_1$  determine a cylindrical decomposition adapted to  $\{P_2\}$ .

- When  $X_1 X_2 \neq 0$ ,  $P_2 = 0$  is equivalent to

$$X_3 = \frac{(X_1 - X_2)^2}{X_1 X_2} \quad \text{or} \quad X_3 = \frac{(X_1 + X_2)^2}{X_1 X_2}.$$

So the zeroes of  $P_2$  are described by two graphs of functions over each open quadrant, and the cylindrical decomposition of  $P_2$  has five cells over each open quadrant. The sign of  $P_2$  in these five cells is  $1, 0, -1, 0, 1$ .

- When  $X_1 = 0, X_2 \neq 0$ , the polynomial  $P_2$  is  $-X_2^4$ . The cylinders over each open half-axis have one cell on which  $P_2$  is negative.
- When  $X_1 \neq 0, X_2 = 0$ , the polynomial  $P_2$  is  $-X_1^4$ . The cylinders over each open half-axis have one cell on which  $P_2$  is negative.
- When  $X_1 = 0, X_2 = 0$ ,  $P_2$  is identically zero. The cylinder over the origin has one cell on which  $P_2$  is zero.

Finally, while the descriptions of the cylindrical decompositions of  $P_2$  and  $P_2$  are identical,  $\text{Zer}(P_1, \mathbb{R}^3)$  and  $\text{Zer}(P_2, \mathbb{R}^3)$  are not homeomorphic:  $\text{Zer}(P_1, \mathbb{R}^3)$  has an open subset which is a semi-algebraic set of dimension 1, and it is not the case for  $\text{Zer}(P_2, \mathbb{R}^3)$ . □

A **semi-algebraic stratification** of a finite family  $S_1, \dots, S_\ell$  of semi-algebraic sets is a finite partition of each  $S_i$  into semi-algebraic sets  $S_{i,j}$  such that

- every  $S_{i,j}$  is a  $\mathcal{S}^\infty$  submanifold,
- the closure of  $S_{i,j}$  in  $S_i$  is the union of  $S_{i,j}$  with some  $S_{i,j'}$ 's where the dimensions of the  $S_{i,j'}$ 's are less than the dimension of  $S_{i,j}$ .

The  $S_{i,j}$  are called **strata** of this stratification. A **cell stratification of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$**  is a stratification of  $\mathbb{R}^k$  for which every stratum  $S_i$  is  $\mathcal{S}^\infty$  diffeomorphic to an open cube  $(0, 1)^{d_i}$  and is also  $\mathcal{P}$ -invariant. A cell stratification of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$  induces a stratification of  $S_1, \dots, S_\ell$  for every finite family  $S_1, \dots, S_\ell$  of  $\mathcal{P}$ -semi-algebraic sets.

**Theorem 5.38.** *For every finite set  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ , there exists a cell stratification of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .*

In Thom’s lemma, the closures of the different “pieces” (points and open intervals) are obtained by relaxing strict inequalities. The key technique to prove Theorem 5.38 is to extend these properties to the case of several variables. In the cylindrical decomposition, when the polynomials are quasi-monic, we can control what happens when we pass from a cylinder  $S \times \mathbb{R}$  to another  $T \times \mathbb{R}$  such that  $T \subset \overline{S}$ . The quasi-monicity is needed to avoid the kind of bad behavior described in Example 5.37.

The following result can be thought of as a multivariate version of Thom’s lemma.

Suppose that

- $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  is closed under differentiation with respect to  $X_k$  and each  $P \in \mathcal{P}$  is quasi-monic with respect to  $X_k$  (see Definition 4.72),
- $S$  and  $S'$  are semi-algebraically connected semi-algebraic subsets of  $\mathbb{R}^{k-1}$ , both  $\text{Elim}_{X_k}(\mathcal{P})$ -invariant, and  $S'$  is contained in the closure of  $S$ .

It follows from Proposition 5.16 that there are continuous semi-algebraic functions  $\xi_1 < \dots < \xi_\ell: S \rightarrow \mathbb{R}$  and  $\xi'_1 < \dots < \xi'_{\ell'}: S' \rightarrow \mathbb{R}$  which describe, for all  $P \in \mathcal{P}$ , the real roots of the polynomials  $P(x, X_k)$  as functions of  $x = (x_1, \dots, x_{k-1})$  in  $S$  or in  $S'$ . Denote by  $\Gamma_j$  and  $\Gamma'_j$  the graphs of  $\xi_j$  and  $\xi'_j$ , respectively. Since  $\mathcal{P}$  is closed under differentiation, there is a polynomial  $P \in \mathcal{P}$  such that, for every  $x \in S$  (resp.  $x \in S'$ ),  $\xi_j(x)$  (resp.  $\xi'_j(x)$ ) is a simple root of  $P(x, X_k)$  for  $P \in \mathcal{P}$  (see Remark 5.35). Denote by  $B_j$  and  $B'_j$  the bands of the cylinders  $S \times \mathbb{R}$  and  $S' \times \mathbb{R}$ , respectively, which are bounded by these graphs.

**Proposition 5.39. [Generalized Thom's Lemma]**

- Every function  $\xi_j$  can be continuously extended to  $S'$ , and this extension coincides with one of the functions  $\xi'_{j'}$ .
- For every function  $\xi'_{j'}$ , there is a function  $\xi_j$  whose extension by continuity to  $S'$  is  $\xi'_{j'}$ .
- For every  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$ , the set

$$\text{Reali}(\sigma, S \times \mathbb{R}) = \{(x, x_k) \in S \times \mathbb{R} \mid \text{sign}(\mathcal{P}(x, x_k)) = \sigma\}$$

is either empty or one of the  $\Gamma_j$  or one of the  $B_j$ . Let  $\text{Reali}(\bar{\sigma}, S \times \mathbb{R})$  be the subset of  $S \times \mathbb{R}$  obtained by relaxing the strict inequalities:

$$\text{Reali}(\bar{\sigma}, S \times \mathbb{R}) = \{(x, x_k) \in S \times \mathbb{R} \mid \text{sign}(\mathcal{P}(x, x_k)) \in \bar{\sigma}\},$$

and let

$$\text{Reali}(\bar{\sigma}; S' \times \mathbb{R}) = \{(x, x_k) \in S' \times \mathbb{R} \mid \text{sign}(\mathcal{P}(x, x_k)) \in \bar{\sigma}\}.$$

If  $\text{Reali}(\sigma, S \times \mathbb{R}) \neq \emptyset$ , we have  $\overline{\text{Reali}(\sigma, S \times \mathbb{R})} \cap (S \times \mathbb{R}) = \text{Reali}(\bar{\sigma}, S \times \mathbb{R})$  and  $\overline{\text{Reali}(\sigma, S \times \mathbb{R})} \cap (S' \times \mathbb{R}) = \text{Reali}(\bar{\sigma}, S' \times \mathbb{R})$ . Moreover,  $\text{Reali}(\sigma, S' \times \mathbb{R})$  is either a graph  $\Gamma'_{j'}$ , or the closure of one of the bands  $B'_{j'}$ , in  $S' \times \mathbb{R}$ .

**Proof:** Let  $x' \in S'$ . Consider one of the functions  $\xi_j$ . Since  $\mathcal{P}$  is closed under differentiation, there is a polynomial  $P \in \mathcal{P}$  such that, for every  $x \in S$ ,  $\xi_j(x)$  is a simple root of

$$P(x, X_k) = a_p X_k^p + a_{p-1}(x) X_k^{p-1} + \dots + a_0(x),$$

(see Remark 5.35). Moreover,  $a_p$  is a non-zero constant. Let

$$M(x') = \max_{i=1, \dots, p} \left( p \left| \frac{a_{p-i}(x')}{a_p} \right| \right)^{1/i}.$$

By Proposition 5.9, and the continuity of  $M$ , there is a neighborhood  $U$  of  $x'$  in  $\mathbb{R}^{k-1}$  such that, for every  $x \in S \cap U$ , we have

$$\xi_j(x) \in [-M(x') - 1, M(x') + 1]$$

Choose a continuous semi-algebraic path  $\gamma$  such that

$$\gamma((0, 1]) \subset S \cap U, \quad \gamma(0) = x'.$$

The semi-algebraic function  $f = \xi_j \circ \gamma$  is bounded and therefore has, by Proposition 3.18, a continuous extension  $\bar{f}$  with

$$\bar{f}(0) \in [-M(x') - 1, M(x') + 1].$$

Let  $\tau_1 = \text{sign}(P'(x, \xi_j(x))), \dots, \tau_p = \text{sign}(P^{(p)}(x, \xi_j(x)))$ , for  $x \in S$  (observe that these signs are constant for  $x \in S$ ). Since every point in the graph of  $\xi_j$  satisfies

$$P(x', x'_k) = 0, \quad \text{sign}(P'(x', x'_k)) = \tau_1, \dots, \quad \text{sign}(P^{(p)}(x', x'_k)) = \tau_d,$$

every point  $(x', x'_k)$  in the closure of the graph of  $\xi_j$  must satisfy

$$P(x', x'_k) = 0, \quad \text{sign}(P'(x', x'_k)) \in \bar{\tau}_1, \dots, \quad \text{sign}(P^{(p)}(x', x'_k)) \in \bar{\tau}_d.$$

By Lemma 5.33 (Thom's lemma), there is at most one  $x'_k$  satisfying these inequalities. Since  $(x', \bar{f}(0))$  is in the closure of the graph of  $\xi_j$ , it follows that  $\xi_j$  extends continuously to  $x'$  with the value  $\bar{f}(0)$ . Hence, it extends continuously to  $S'$ , and this extension coincides with one of the functions  $\xi_{j'}$ . This proves the first item.

We now prove the second item. Choose a function  $\xi_{j'}$ . Since  $\xi_{j'}$  is a simple root of some polynomial  $P$  in the set, by Proposition 3.10 there is a function  $\xi_j$ , also a root of  $P$ , whose continuous extension to  $S'$  is  $\xi_{j'}$ .

We now turn to the third item. The properties of  $\text{Reali}(\sigma, S \times \mathbb{R})$  and  $\text{Reali}(\bar{\sigma}, S \times \mathbb{R})$  are straightforward consequences of Thom's lemma, since  $P \in \mathcal{P}$  has constant sign on each graph  $\Gamma_j$  and each band  $B_j$ . The closure of  $B_j$  in  $S \times \mathbb{R}$  is  $\Gamma_j \cup B_j \cup \Gamma_{j+1}$ , where  $\Gamma_0 = \Gamma_{\ell+1} = \emptyset$  and therefore it is obvious that  $\overline{\text{Reali}(\sigma, S \times \mathbb{R})} \cap (S' \times \mathbb{R}) \subset \text{Reali}(\bar{\sigma}, S' \times \mathbb{R})$ . It follows from 1 and 2 that  $\overline{\text{Reali}(\sigma, S \times \mathbb{R})} \cap (S' \times \mathbb{R})$  is either a graph  $\Gamma'_j$  or the closure of one of the bands  $B'_{j'}$  in  $S' \times \mathbb{R}$ .

By Thom's lemma, this is also the case for  $\text{Reali}(\bar{\sigma}, S' \times \mathbb{R})$ . It remains to check that it cannot happen that  $\text{Reali}(\bar{\sigma}, S' \times \mathbb{R})$  is the closure of a band  $B'_{j'}$ , and  $\overline{\text{Reali}(\sigma, S \times \mathbb{R})} \cap (S' \times \mathbb{R})$  is one of the graphs  $\Gamma'_j$  or  $\Gamma'_{j'+1}$ . In this case, all  $\sigma(P)$  should be different from zero and the sign of  $P$  should be  $\sigma(P)$  on every sufficiently small neighborhood  $V$  of a point  $x'$  of  $B'_{j'}$ . This implies that  $V \cap (S \times \mathbb{R}) \subset R(\sigma, S \times \mathbb{R})$  and, hence,  $x' \in \overline{R(\sigma, S \times \mathbb{R})}$ , which is impossible. □

**Proposition 5.40.** *Let  $\mathcal{P}^* = \bigcup_{i=1}^k \mathcal{P}_i$  be finite sets of non-zero polynomials such that:*

- $\mathcal{P}_k$  contains  $\mathcal{P}$ ,
- or each  $i$ ,  $\mathcal{P}_i$  is a subset of  $\mathbb{R}[X_1, \dots, X_i]$  that is closed under differentiation and quasi-monic with respect to  $X_i$ ,



– for  $i \leq k$ ,  $\text{Elim}_{X_i}(\mathcal{P}_i) \subset \mathcal{P}_{i-1}$ .

Writing  $\mathcal{P}_{\leq i} = \bigcup_{j \leq i} \mathcal{P}_j$ , the families  $\mathcal{S}_i$ , for  $i = 1, \dots, k$ , consisting of all  $\text{Reali}(\sigma)$  with  $\sigma$  a realizable sign conditions on  $\mathcal{P}_{\leq i}$  constitute a cylindrical decomposition of  $\mathbb{R}^k$  that is a cell stratification of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .

**Proof:** The proof of the proposition is a simple induction on  $k$ . The preceding Proposition 5.39 (Generalized Thom's Lemma) provides the induction step and Thom's lemma the base case for  $k = 1$ . To show that the dimension condition is satisfied, observe that if  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$  and  $\text{Reali}(\sigma) \neq \emptyset$ , then  $\text{Reali}(\bar{\sigma})$  is the union of  $\text{Reali}(\sigma)$  and some  $\text{Reali}(\sigma')$ ,  $\sigma' \neq \sigma$ .

Since  $\text{Reali}(\sigma)$  (resp.  $\text{Reali}(\sigma')$ ) is a cell of a cylindrical decomposition,  $\text{Reali}(\sigma)$  (resp.  $\text{Reali}(\sigma')$ ) is semi-algebraically homeomorphic to  $(0, 1)^{d(\sigma)}$  (resp.  $(0, 1)^{d(\sigma')}$ ). That  $d(\sigma') < d(\sigma)$  is easily seen by induction.  $\square$

A family  $\mathcal{P}$  is a **stratifying family** if it satisfies the hypothesis of Proposition 5.40.

The theorem above holds for a stratifying family of polynomials. But we shall now see that it is always possible to convert a finite set of polynomials to a quasi-monic set by making a suitable linear change of coordinates. By successively converting to quasi-monic polynomials, closing under differentiation and applying Elim, we arrive at a stratifying family.

**Proposition 5.41.** *Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ . There is a linear automorphism  $u: \mathbb{R}^k \rightarrow \mathbb{R}^k$  and a stratifying family of polynomials  $\mathcal{Q}^* = \cup_{i=1, \dots, k} \mathcal{Q}_i$  such that  $P(u(X)) \in \mathcal{Q}_k$  for all  $P \in \mathcal{P}$  (where  $X = (X_1, \dots, X_k)$ ).*

**Proof:** By Lemma 4.73, there is a linear change of variables  $v$  such that, for all  $P \in \mathcal{P}$ , the polynomial  $P(v(X))$  is quasi-monic with respect to  $X_k$ .

Let  $\mathcal{Q}_k$  consist of all polynomials  $P(v(X))$  for  $P \in \mathcal{P}$  together with all their non-zero derivatives of every order with respect to  $X_k$ . Using induction, applied to  $\text{Elim}_{X_k}(\mathcal{Q}_k)$ , there is a linear automorphism  $u': \mathbb{R}^{k-1} \rightarrow \mathbb{R}^{k-1}$  and a stratifying family of polynomials  $\cup_{1 \leq i \leq k-1} \mathcal{R}_i$  such that  $Q(u'(X')) \in \mathcal{R}_{k-1}$  for every  $Q \in \text{Elim}_{X_k}(\mathcal{Q}_k)$ , where  $X' = (X_1, \dots, X_{k-1})$ . Finally, set  $u = (u' \times \text{Id}) \circ v$  (where  $u' \times \text{Id}(X', X_k) = (u'(X'), X_k)$ ),  $\mathcal{Q}_j = \{R(X) \mid R \in \mathcal{R}_j\}$  for  $j \leq k-1$ .  $\square$

We are now ready for the proof of Theorem 5.38.

**Proof of Theorem 5.38:** Use Proposition 5.41 to get a linear automorphism  $u: \mathbb{R}^k \rightarrow \mathbb{R}^k$  and a stratifying family  $\mathcal{Q}^*$  that contains

$$\mathcal{Q} = \{P(u(X)) \mid P \in \mathcal{P}\}$$

in order to obtain, by Proposition 5.40, a cell stratification adapted to  $\mathcal{Q}$ . Clearly,  $u^{-1}$  converts this cell stratification to one adapted to  $\mathcal{P}$ .  $\square$

Theorem 5.38 has consequences for the dimension of semi-algebraic sets.

**Theorem 5.42.** *Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set. Then,*

$$\begin{aligned} \dim(\overline{S}) &= \dim(S), \\ \dim(\overline{S} \setminus S) &< \dim(S). \end{aligned}$$

**Proof:** This is clear from Proposition 5.28 and Theorem 5.38, since the closure of a stratum is the union of this stratum and of a finite number of strata of smaller dimensions.  $\square$

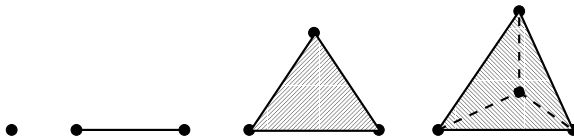
### 5.6 Simplicial Complexes

We first recall some basic definitions and notations about simplicial complexes.

Let  $a_0, \dots, a_d$  be points of  $\mathbb{R}^k$  that are **affinely independent** (which means that they are not contained in any affine subspace of dimension  $d - 1$ ). The  **$d$ -simplex** with vertices  $a_0, \dots, a_d$  is

$$[a_0, \dots, a_d] = \left\{ \lambda_0 a_0 + \dots + \lambda_d a_d \mid \sum_{i=0}^d \lambda_i = 1 \text{ and } \lambda_0, \dots, \lambda_d \geq 0 \right\}.$$

Note that the dimension of  $[a_0, \dots, a_d]$  is  $d$ .



**Fig. 5.2.** Zero, one, two, and three dimensional simplices

An  $e$ -**face** of the  $d$ -simplex  $s = [a_0, \dots, a_d]$  is any simplex  $s' = [b_0, \dots, b_e]$  such that

$$\{b_0, \dots, b_e\} \subset \{a_0, \dots, a_d\}.$$

The face  $s'$  is a **proper face** of  $s$  if  $\{b_0, \dots, b_e\} \neq \{a_0, \dots, a_d\}$ . The 0-faces of a simplex are its **vertices**, the 1-faces are its **edges**, and the  $(d - 1)$ -faces of a  $d$ -simplex are its **facets**. We also include the empty set as a simplex of dimension  $-1$ , which is a face of every simplex. If  $s'$  is a face of  $s$  we write  $s' \prec s$ .

The open simplex corresponding to a simplex  $s$  is denoted  $s^\circ$  and consists of all points of  $s$  which do not belong to any proper face of  $s$ :

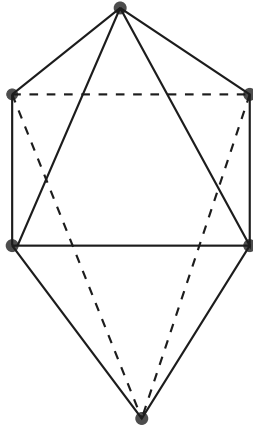
$$s^\circ = (a_0, \dots, a_d) = \left\{ \lambda_0 a_0 + \dots + \lambda_d a_d \mid \sum_{i=0}^d \lambda_i = 1 \text{ and } \lambda_0 > 0, \dots, \lambda_d > 0 \right\}.$$

which is the interior of  $[a_0, \dots, a_d]$ . By convention, if  $s$  is a 0 – simplex then  $s^\circ = s$ .

The **barycenter** of a  $d$  – simplex  $s = [a_0, \dots, a_d]$  in  $\mathbb{R}^k$  is the point  $\text{ba}(s) \in \mathbb{R}^k$  defined by  $\text{ba}(s) = 1/(d+1) \sum_{0 \leq i \leq d} a_i$ .

A **simplicial complex**  $K$  in  $\mathbb{R}^k$  is a finite set of simplices in  $\mathbb{R}^k$  such that  $s, s' \in K$  implies

- every face of  $s$  is in  $K$ ,
- $s \cap s'$  is a common face of both  $s$  and  $s'$ .



**Fig. 5.3.** A two dimensional simplicial complex homeomorphic to  $S^2$

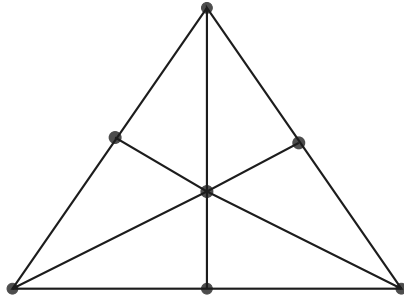
The set  $|K| = \bigcup_{s \in K} s$ , which is clearly a semi-algebraic subset of  $\mathbb{R}^k$ , is called the the **realization of  $K$** . Note that the realization of  $K$  is the disjoint union of its open simplices. A **polyhedron** in  $\mathbb{R}^k$  is a subset  $P$  of  $\mathbb{R}^k$  such that there exists a simplicial complex  $K$  in  $\mathbb{R}^k$  with  $P = |K|$ . Such a  $K$  is called a **simplicial decomposition** of  $P$ .

Let  $K$  and  $L$  be two simplicial complexes. Then  $L$  is called a **subdivision** of  $K$  if

- $|L| = |K|$ ,
- for every simplex  $s \in L$  there is a simplex  $s' \in K$  such that  $s \subset s'$ .

Given a simplicial complex  $K$ , an **ascending sequence of simplices** is a collection of simplices  $\{s_0, s_1, \dots, s_j\}$  such that  $s_0 \prec s_1 \prec \dots \prec s_j$ .

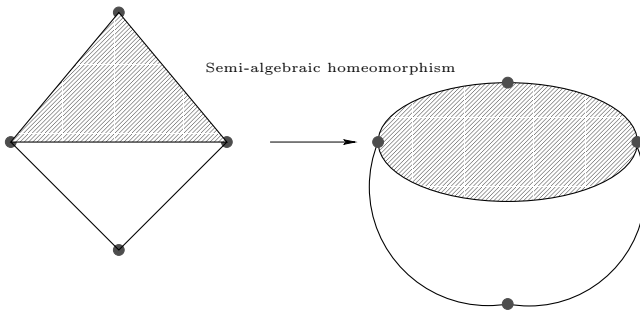
Let  $K$  be a simplicial complex. Let  $\text{ba}(K)$  denote the set of simplices that are spanned by the barycenters of some ascending sequence of simplices of  $K$ . Thus for every ascending sequence of simplices in  $K$ ,  $s_0 \prec s_1 \prec \dots \prec s_j$ , we include in  $K'$  the simplex  $[\text{ba}(s_0), \dots, \text{ba}(s_j)]$ , and we call  $\text{ba}(s_j)$  the **leading vertex** of  $[\text{ba}(s_0), \dots, \text{ba}(s_j)]$ . It is easy to check that  $\text{ba}(K)$  is a simplicial complex, called the **barycentric subdivision** of  $K$ .



**Fig. 5.4.** The barycentric subdivision of a two dimensional simplex

## 5.7 Triangulation

A **triangulation** of a semi-algebraic set  $S$  is a simplicial complex  $K$  together with a semi-algebraic homeomorphism  $h$  from  $|K|$  to  $S$ . We next prove that any closed and bounded semi-algebraic set can be triangulated. In fact, we prove a little more, which will be useful for technical reasons. The triangulation will also be a stratification of  $S$  which respects any given finite collection of semi-algebraic subsets of  $S$ , i.e. the images of the open simplices will be the strata and each of the specified subsets of  $S$  will be stratified as well.



**Fig. 5.5.** Semi-algebraic triangulation

A triangulation of  $S$  **respecting a finite family of semi-algebraic sets**  $S_1, \dots, S_q$  contained in  $S$  is a triangulation  $K, h$  such that each  $S_j$  is the union of images by  $h$  of open simplices of  $K$ .

**Theorem 5.43. [Triangulation]** *Let  $S \subset \mathbb{R}^k$  be a closed and bounded semi-algebraic set, and let  $S_1, \dots, S_q$  be semi-algebraic subsets of  $S$ . There exists a triangulation of  $S$  respecting  $S_1, \dots, S_q$ . Moreover, the vertices of  $K$  can be chosen with rational coordinates.*

**Proof:** We first prove the first part of the statement. The proof is by induction on  $k$ . For  $k = 1$ , let  $|K| = S$ , taking as open simplices the points and bounded open intervals which constitute  $S$ .

We prove the result for  $k > 1$  supposing that it is true for  $k - 1$ . After a linear change of variables as in Proposition 5.41, we can suppose that  $S$  and the  $S_j$  are the union of strata of a stratifying set of polynomials  $\mathcal{P}$ . Thus  $\mathbb{R}^{k-1}$  can be decomposed into a finite number of semi-algebraically connected semi-algebraic sets  $C_i$ , and there are semi-algebraic and continuous functions  $\xi_{i,1} < \dots < \xi_{i,\ell_i}: C_i \rightarrow \mathbb{R}$  describing the roots of the non-zero polynomials among  $P(x, X_k)$ ,  $P \in \mathcal{P}_k$ , as functions of  $x \in C_i$ . We know that  $S$ , and the  $S_j$ , are unions of some graphs of  $\xi_{i,j}$  and of some bands of cylinders  $C_i \times \mathbb{R}$  between these graphs. Denote by  $\pi: \mathbb{R}^k \rightarrow \mathbb{R}^{k-1}$  the projection which forgets the last coordinate. The set  $\pi(S)$  is closed and bounded, semi-algebraic, and the union of some  $C_i$ ; also, each  $\pi(S_j)$  is the union of some  $C_i$ . By the induction hypothesis, there is a triangulation  $g: |L| \rightarrow \pi(S)$  (where  $g$  is a semi-algebraic homeomorphism,  $L$  a simplicial complex in  $\mathbb{R}^{k-1}$ ) such that each  $C_i \subset \pi(S)$  is a union of images by  $g$  of open simplices of  $L$ . Thus, at the top level,  $\mathbb{R}^k$  is decomposed into cylinders over sets of the form  $g(t^\circ)$  for  $t$  a simplex of  $L$ .

We next extend the triangulation of  $\pi(S)$  to one for  $S$ . For every  $t$  in  $L$  we construct a simplicial complex  $K_t$  and a semi-algebraic homeomorphism

$$h_t: |K_t| \longrightarrow \overline{S \cap (g(t^\circ) \times \mathbb{R})}.$$

Fix a  $t$  in  $L$ , say  $t = [b_0, \dots, b_d]$ , and let  $\xi: g(t^\circ) \rightarrow \mathbb{R}$  be a function of the cylindrical decomposition whose graph is contained in  $S$ . We are in the situation of Proposition 5.39, and we know that  $\xi$  can be extended continuously to  $\bar{\xi}$  defined on the closure of  $g(t^\circ)$  which is  $g(t)$ . Define  $a_i = (b_i, \bar{\xi}(g(b_i))) \in \mathbb{R}^k$  for  $i = 0, \dots, d$ , and let  $s_\xi$  be the simplex  $[a_0, \dots, a_d] \subset \mathbb{R}^k$ . The simplex  $s_\xi$  will be a simplex of the complex  $K_t$  we are constructing. Define  $h_t$  on  $s_\xi$  by

$$h_t(\lambda_0 a_0 + \dots + \lambda_d a_d) = (y, \bar{\xi}(y)), \quad \text{where } y = g(\lambda_0 b_0 + \dots + \lambda_d b_d).$$

If  $\xi': g(t^\circ) \rightarrow \mathbb{R}$  is another function of the cylindrical decomposition whose graph is contained in  $S$ , define  $s_{\xi'} = [a'_0, \dots, a'_d]$  in the same way. It is important that  $s_{\xi'}$  not coincide with  $s_\xi$ . At least one of the  $a'_i$  must differ from the corresponding  $a_i$ . Similarly when the restrictions of  $\bar{\xi}$  and  $\bar{\xi}'$  to a face  $r$  of  $t$  are not the same, we require that on at least one vertex  $b_i$  of  $r$ , the values of  $\bar{\xi}$  and  $\bar{\xi}'$  are different (so that the corresponding  $a_i$  and  $a'_i$  are distinct). Thus we require that on every simplex  $t$  of  $L$ , if  $\xi$  and  $\xi'$  are two distinct functions  $g(t^\circ) \rightarrow \mathbb{R}$  of the cylindrical decomposition then there exists a vertex  $b$  of  $t$  such that  $\bar{\xi}(g(b)) \neq \bar{\xi}'(g(b))$ . It is clear that this requirement will be satisfied if we replace  $L$  by its barycentric division  $\text{ba}(L)$ . Hence, after possibly making this replacement, we can assume that our requirement is satisfied by  $L$ .

Now consider a band between two graphs of successive functions

$$\xi < \xi': g(t^\circ) \longrightarrow \mathbb{R}$$

contained in  $S$  (note that an unbounded band cannot be contained in  $S$ , since  $S$  is closed and bounded). Let  $P$  be the polyhedron above  $t$  whose bottom face is  $s_\xi$  and whose top face is  $s_{\xi'}$ . This polyhedron  $P$  has a simplicial decomposition

$$P = \bigcup_{i=0}^d [a'_0, \dots, a'_i, a_i, \dots, a_d].$$

Note that it may happen that  $a'_i = a_i$  in which case we understand

$$[a'_0, \dots, a'_i, a_i, \dots, a_d]$$

to be the  $d - 1$ -simplex

$$[a'_0, \dots, a'_i, a_{i+1}, \dots, a_d].$$

The complex  $K_t$  we are constructing contains the simplices (and their faces) of this simplicial decomposition of  $P$ . We define  $h_t$  on  $P$  by the condition that the segment  $[\lambda_0 a_0 + \dots + \lambda_d a_d, \lambda_0 a'_0 + \dots + \lambda_d a'_d]$  is sent by an affine function to  $[(y, \bar{\xi}(y)), (y, \bar{\xi}'(y))]$ , where

$$y = g(\lambda_0 b_0 + \dots + \lambda_d b_d).$$

Having constructed  $K_t$  and  $h_t$  for each simplex  $t$  of  $L$ , it remains to prove that these  $K_t$  and  $h_t$  can be glued together to give  $K$  and  $h$  as a triangulation of  $S$ . We next show that it is possible if we first choose a total order for all vertices of  $L$  and then label the simplices of  $L$  compatibly with this total order.

It is enough to check this for a simplex  $t$  and one of its faces  $r$ . The first thing to notice is that if we have a simplex  $s_\eta$  in  $K_r$  that is sent by  $h_r$  onto the closure of the graph  $\eta: g(r^\circ) \rightarrow \mathbb{R}$ , a function of the algebraic decomposition, and if  $s_\eta$  meets  $|K_t|$ , then it is a simplex of  $K_t$ : indeed in this case  $\eta$  coincides with one of the  $\bar{\xi}$  on  $g(r^\circ)$  by point 2 of Proposition 5.39 (for  $\xi: g(t^\circ) \rightarrow \mathbb{R}$ ,  $s_\xi$  simplex of  $K_t$  and  $s_\eta$  a facet of  $s_\xi$ ). For this reason, it is also the case that  $h_t$  and  $h_r$  coincide on  $|K_t| \cap |K_r|$ . What remains to verify is that the simplicial complex of the polyhedron  $P$  in  $t \times \mathbb{R}$  (see above) induces the simplicial decomposition of the polyhedron  $P \cap (r \times \mathbb{R})$ . This is the case if the simplicial decomposition  $P = \bigcup_{i=0}^d [a'_0, \dots, a'_i, a_i, \dots, a_d]$  is compatible with a fixed total order on the vertices of  $L$ .

It remains to prove that there exists a simplicial complex  $L$  with rational coordinates such that  $|K|$  and  $|L|$  are semi-algebraically homeomorphic. The proof is by induction on  $k$ . When  $k = 1$ , the semi-algebraic subsets  $S, S_1, \dots, S_q$  are a finite number of points and intervals and the claim is clear. The inductive steps uses the cylindrical structure of the constructed triangulation. □

The following corollary of Theorem 5.43 will be used in the proof of Theorem 5.46.

**Proposition 5.44.** *Let  $S \subset \mathbb{R}^k$  be a closed and bounded semi-algebraic set, and let  $S_1, \dots, S_q$  be semi-algebraic subsets of  $S$ , such that  $S$  and the  $S_j$  are given by boolean combinations of sign conditions over polynomials that are all either monic in the variable  $X_k$  or independent of  $X_k$ . Let  $\pi$  be the projection of  $\mathbb{R}^k$  to  $\mathbb{R}^{k-1}$  that forgets the last factor. There are semi-algebraic triangulations*

$$\Phi: |K| = \bigcup_{p=1}^{s'} |s_p| \rightarrow S, \quad |K| \subset \mathbb{R}^k$$

and

$$\Psi: |L| = \bigcup_{\ell=1}^s |t_\ell| \rightarrow \pi(S), \quad |L| \subset \mathbb{R}^{k-1}$$

such that  $\pi \circ \Phi(x) = \Psi \circ \pi(x)$  for  $x \in |K|$ , and each  $S_j$  is the union of some  $\Phi(s_i^?)$ .

### 5.8 Hardt’s Triviality Theorem and Consequences

Hardt’s triviality theorem is the following.

**Theorem 5.45. [Hardt’s triviality theorem]** *Let  $S \subset \mathbb{R}^m$  and  $T \subset \mathbb{R}^k$  be semi-algebraic sets. Given a continuous semi-algebraic function  $f: S \rightarrow T$ , there exists a finite partition of  $T$  into semi-algebraic sets  $T = \bigcup_{i=1}^r T_i$ , so that for each  $i$  and any  $x_i \in T_i$ ,  $T_i \times f^{-1}(x_i)$  is semi-algebraically homeomorphic to  $f^{-1}(T_i)$ .*

For technical reasons, we prove the slightly more general:

**Theorem 5.46. [Semi-algebraic triviality]** *Let  $S \subset \mathbb{R}^m$  and  $T \subset \mathbb{R}^k$  be semi-algebraic sets. Given a continuous semi-algebraic function  $f: S \rightarrow T$  and  $S_1, \dots, S_q$  semi-algebraic subsets of  $S$ , there exists a finite partition of  $T$  into semi-algebraic sets  $T = \bigcup_{i=1}^r T_i$ , so that for each  $i$  and any  $x_i \in T_i$ ,  $T_i \times f^{-1}(x_i)$  is semi-algebraically homeomorphic to  $f^{-1}(T_i)$ . More precisely, writing  $F_i = f^{-1}(x_i)$ , there exist semi-algebraic subsets  $F_{i,1}, \dots, F_{i,q}$  of  $F_i$  and a semi-algebraic homeomorphism  $\theta_i: T_i \times F_i \rightarrow f^{-1}(T_i)$  such that  $f \circ \theta_i$  is the projection mapping  $T_i \times F_i \rightarrow T_i$  and such that*

$$\theta_i(T_i \times F_{i,j}) = S_j \cap f^{-1}(T_i).$$

**Proof:** We may assume without loss of generality that

- $S$  and  $T$  are both bounded (using if needed homeomorphisms of the form  $x \mapsto x/(1 + \|x\|)$ , which are obviously semi-algebraic),
- $S$  is a semi-algebraic subset of  $\mathbb{R}^{m+k}$  and  $f$  is the restriction to  $S$  of the projection mapping  $\Pi: \mathbb{R}^{m+k} \rightarrow \mathbb{R}^k$  that forgets the first  $m$  coordinates, (replacing  $S$  by the graph of  $f$  which is semi-algebraically homeomorphic to  $S$ ).

The proof proceeds by induction on the lexicographic ordering of the pairs  $(m, k)$ .

The sets  $S$  and  $S_j$  are given by boolean combinations of sign conditions over a finite number of polynomials  $\mathcal{P} \subset \mathbb{R}[X, Y]$ , where  $X = (X_1, \dots, X_m)$  and  $Y = (Y_1, \dots, Y_k)$ . Making, if needed, a linear substitution of the variables of the  $Y$ 's only as in Proposition 5.41, one may suppose that each  $P \in \mathcal{P}$  can be written

$$g_{P,0}(X) Y_k^{d(P)} + g_{P,1}(X, Y') Y_k^{d(P)-1} + \dots + g_{P,d(P)}(X, Y'),$$

where  $Y' = (Y_1, \dots, Y_{k-1})$ , with  $g_{P,0}(X)$  not identically zero. Let

$$A(X) = \prod_{P \in \mathcal{P}} g_{P,0}(X).$$

The dimension of the semi-algebraic set  $T'' = \{x \in T \mid A(x) = 0\}$  is strictly smaller than  $m$ . By Theorem 5.19, this set can be written as the finite union of sets of the form  $\varphi((0, 1)^d)$  where  $\varphi$  is a semi-algebraic homeomorphism and  $d < m$ . Taking the inverse image under  $\varphi$ , we have to deal with a subset of  $\mathbb{R}^d$ , and our induction hypothesis takes care of this case.

It remains to handle what happens above  $T' = T \setminus T''$ . We multiply each polynomial in  $\mathcal{P}$  by a convenient product of powers of  $g_{Q,0}(X)$ ,  $Q \in \mathcal{P}$ , so that the leading coefficient of  $P$  becomes  $(A(X)Y_k)^{d(P)}$ . Replacing  $A(X)Y_k$  by  $Z_k$  defines a semi-algebraic homeomorphism from  $S \cap (T' \times \mathbb{R}^k)$  onto a bounded semi-algebraic set  $S' \subset \mathbb{R}^{m+k}$ . Denote by  $S'_j$  the image of  $S_j \cap (T' \times \mathbb{R}^k)$  under this homeomorphism. Now, the sets  $S'$  and  $S'_j$  are both given by boolean combinations of sign conditions over polynomials that are all either quasi-monic in the variable  $Z_k$  or independent of  $Z_k$ . Up to a linear substitution of the variables involving only the variables  $X$  and  $(Y_1, \dots, Y_{k-1})$ , one may suppose that  $S'$  and the  $S'_j$  are given by boolean combinations of sign conditions over polynomials of a stratifying family. By Proposition 5.40,  $\overline{S'}$  is also given by a boolean combination of sign conditions over the same polynomials.

One can now apply Corollary 5.45 to  $\overline{S'}$  and  $S'_0 = S', S'_1, \dots, S'_q$ : there are semi-algebraic triangulations

$$\Phi: |K| = \bigcup_{p=1}^{s'} |s_p| \rightarrow \overline{S'}, \quad |K| \subset \mathbb{R}^{m+k}$$

and

$$\Psi: |L| = \bigcup_{\ell=1}^s |t_\ell| \rightarrow \pi(\overline{S'}), \quad |L| \subset \mathbb{R}^{m+k-1}$$

such that  $\pi \circ \Phi(x) = \Psi \circ \pi(x)$  if  $x \in |K|$ , and each  $S'_j$  ( $j = 0, \dots, q$ ) is the union of some  $\Phi(s_i^c)$ .



We now apply the induction hypothesis to  $\pi(\overline{S'})$ , with the subsets  $\Psi(t_\ell^\circ)$  and the projection mapping  $\Pi': \mathbb{R}^{m-1+k} \rightarrow \mathbb{R}^k$ . We obtain a finite partition of  $\mathbb{R}^k$  into semi-algebraic sets  $(T'_i)_{i=1, \dots, r}$ . We also obtain semi-algebraic sets  $G_i, G_{i,0}, G_{i,1}, \dots, G_{i,s}$  with  $G_{i,\ell} \subset G_i \subset \mathbb{R}^{m-1}$  and semi-algebraic homeomorphisms  $\rho_i: T'_i \times G_i \rightarrow \Pi'^{-1}(T'_i) \cap \pi(\overline{S'})$  such that  $\Pi' \circ \rho_i$  is the projection mapping  $T'_i \times G_i \rightarrow T'_i$ . Moreover, for every  $\ell$ ,  $\rho_\ell(T'_i \times G_{i,\ell}) = \Pi'^{-1}(T'_i) \cap \Psi(t_\ell^\circ)$ . Let us fix  $i$ , and let  $x_1$  be a point of  $T'_i$ . One may suppose that

$$G_i = \Pi'^{-1}(x_1) \cap \pi(\overline{S'})$$

and that if  $(x_1, y') \in G_i$ , then  $\rho_i(x_1, (x_1, y')) = (x_1, y')$ . Let us then set  $F'_i = \Pi^{-1}(x_1) \cap S'$  and  $F'_{i,j} = \Pi^{-1}(x_1) \cap S'_j$ . It remains to build

$$\theta_i: T'_i \times F'_i \rightarrow \Pi^{-1}(T'_i) \cap \overline{S'}$$

Let  $x \in T'_i$  and  $(x_1, y') \in G_i$ ;  $(x_1, y')$  belongs to one of the  $\Psi(t_\ell^\circ)$ , say  $\Psi(t_1^\circ)$ , and  $\rho_i(x, (x_1, y')) \in \Psi(t_1^\circ)$ . By the properties of the triangulations  $\Phi$  and  $\Psi$ , the intersections with the  $\Phi(s_p)$  decompose

$$\pi^{-1}(x_1, y') \cap S' \text{ and } \pi^{-1}(\rho_\ell(x, (x_1, y'))) \cap \overline{S'}$$

in the same way:  $\theta_i$  maps affinely the segment

$$\{x\} \times (\pi^{-1}(x_1, y') \cap \Phi(s_p)) \subset T'_i \times F'_i$$

(which is possibly either a point or empty) onto the segment

$$\pi^{-1}(\rho_i(x, (x_1, y'))) \cap \Phi(s_p).$$

We leave it to the reader to verify that the  $\theta_i$  built in this way is a semi-algebraic homeomorphism and that  $\theta_i(T'_i \times F'_{i,j}) = \Pi^{-1}(T'_i) \cap S'_j$ .  $\square$

Theorem 5.46 (Semi-algebraic triviality) makes it possible to give an easy proof that the number of topological types of algebraic subsets of  $\mathbb{R}^k$  is finite if one fixes the maximum degree of the polynomials.

**Theorem 5.47. [Finite topological types]** *Let  $k$  and  $d$  be two positive integers. Let  $\mathcal{M}(k, d)$  be the set of algebraic subsets  $V \subset \mathbb{R}^k$  such that there exists a finite set  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  with  $V = \text{Zer}(\mathcal{P}, \mathbb{R}^k)$  and  $\deg(P) \leq d$  for every  $P \in \mathcal{P}$ . There exist a finite number of algebraic subsets  $V_1, \dots, V_s$  of  $\mathbb{R}^k$  in  $\mathcal{M}(k, d)$  such that for every  $V$  in  $\mathcal{M}(k, d)$  there exist  $i$ ,  $1 \leq i \leq s$ , and a semi-algebraic homeomorphism  $\varphi: \mathbb{R}^k \rightarrow \mathbb{R}^k$  with  $\varphi(V_i) = V$ .*

**Proof:** The set  $\mathcal{M}(k, d)$  is contained in the set  $\mathcal{F}$  of algebraic subsets of  $\mathbb{R}^k$  given by a single equation of degree  $\leq 2d$  because

$$\text{Zer}(\mathcal{P}, \mathbb{R}^k) = \text{Zer}\left(\sum_{P \in \mathcal{P}} P^2, \mathbb{R}^k\right).$$

One parametrizes the set  $\mathcal{F}$  by the space  $\mathbb{R}^N$  of coefficients of the equation: abusing notation,  $P$  denotes the point of  $\mathbb{R}^N$  whose coordinates are the coefficients of  $P$ . Let  $S = \{(P, x) \in \mathbb{R}^N \times \mathbb{R}^k \mid P(x) = 0\}$ . The set  $S$  is algebraic. Let  $\Pi: \mathbb{R}^N \times \mathbb{R}^k \rightarrow \mathbb{R}^N$  be the canonical projection mapping. One has  $\Pi^{-1}(P) \cap S = \{P\} \times \text{Zer}(P, \mathbb{R}^k)$ . Theorem 5.46 applied to the projection mapping  $\Pi: \mathbb{R}^N \times \mathbb{R}^k \rightarrow \mathbb{R}^N$  and to the subset  $S$  of  $\mathbb{R}^N \times \mathbb{R}^k$  gives the result.  $\square$

Another consequence of Theorem 5.46 (Semi-algebraic triviality) is the theorem of local conic structure of the semi-algebraic sets.

**Theorem 5.48. [Local conic structure]** *Let  $E$  be a semi-algebraic subset of  $\mathbb{R}^k$  and  $x$  a non-isolated point of  $E$ . Then there exist  $r \in \mathbb{R}$ ,  $r > 0$ , and for every  $r', 0 < r' \leq r$ , a semi-algebraic homeomorphism  $\varphi: \overline{B}_k(x, r') \rightarrow \overline{B}_k(x, r')$  such that:*

- $\|\varphi(y) - x\| = \|y - x\|$  for every  $y \in \overline{B}_k(x, r')$ ,
- $\varphi|_{S^{k-1}(x, r')}$  is the identity mapping,
- $\varphi(E \cap \overline{B}_k(x, r'))$  is a cone with vertex  $x$  and base  $E \cap S^{k-1}(x, r')$ .

**Proof:** Apply Theorem 5.46 (Semi-algebraic triviality) with  $S = \mathbb{R}^k$ ,  $S_1 = E$ , and  $f: S \rightarrow \mathbb{R}$  defined by  $f(y) = \|y - x\|$  to deduce that there exists  $r > 0$  and for every  $r', 0 < r' \leq r$ , a semi-algebraic homeomorphism

$$\theta: (0, r'] \times S^{k-1}(x, r') \rightarrow \overline{B}_k(x, r') \setminus \{x\}$$

such that, for every  $y$  in  $S^{k-1}(x, r')$ ,  $\|\theta(t, y) - x\| = t$  for  $t \in (0, r']$ ,  $\theta(r', y) = y$ , and  $\theta((0, r'] \times (E \cap S^{k-1}(x, r'))) = E \cap \overline{B}_k(x, r) \setminus \{x\}$ . It is then easy to build  $\varphi$ .  $\square$

Let  $S$  be a closed semi-algebraic set and  $T$  a closed semi-algebraic subset of  $S$ . A **semi-algebraic deformation retraction** from  $S$  to  $T$ , is a continuous semi-algebraic function  $h: S \times [0, 1] \rightarrow S$  such that  $h(-, 0)$  is the identity mapping of  $S$ , such that  $h(-, 1)$  has its values in  $T$  and such that for every  $t \in [0, 1]$  and every  $x$  in  $T$ ,  $h(x, t) = x$ .

**Proposition 5.49. [Conic structure at infinity]** *Let  $S$  be a closed semi-algebraic subset of  $\mathbb{R}^k$ . There exists  $r \in \mathbb{R}$ ,  $r > 0$ , such that for every  $r', r' \geq r$ , there is a semi-algebraic deformation retraction from  $S$  to  $S_{r'} = S \cap \overline{B}_k(0, r')$  and a semi-algebraic deformation retraction from  $S_{r'}$  to  $S_r$ .*

**Proof:** Let us suppose that  $S$  is not bounded. Through an inversion mapping  $\varphi: \mathbb{R}^k \setminus \{0\} \rightarrow \mathbb{R}^k \setminus \{0\}$ ,  $\varphi(x) = x/\|x\|^2$ , which is obviously semi-algebraic, one can reduce to the property of local conic structure for  $\varphi(S) \cup \{0\}$  at 0.  $\square$

**Proposition 5.50.** *Let  $f: S \rightarrow T$  be a semi-algebraic function that is a local homeomorphism. There exists a finite cover  $S = \bigcup_{i=1}^n U_i$  of  $S$  by semi-algebraic sets  $U_i$  that are open in  $S$  and such that  $f|_{U_i}$  is a homeomorphism for every  $i$ .*

**Proof:** We assume, as in the proof of Theorem 5.46, that  $T$  is bounded and that the partition  $T = \bigcup_{\ell=1}^r T_\ell$  is induced by a semi-algebraic triangulation  $\Phi: |K| = \bigcup_{\ell=1}^s |s_\ell| \rightarrow \bar{T}$  such that  $T_\ell = \Phi(s_\ell^0)$ . We then replace  $T$  by  $Z = \bigcup_{\ell=1}^r |s_\ell^0|$  and set  $g = \Phi^{-1} \circ f$ . There are semi-algebraic homeomorphisms  $\theta_\ell: s_\ell^0 \times F_\ell \rightarrow g^{-1}(s_\ell^0)$  such that  $g \circ \theta_\ell$  is the projection mapping  $s_\ell^0 \times F_\ell \rightarrow s_\ell^0$  by Theorem 5.46 (Semi-algebraic triviality). Each  $F_\ell$  consists of a finite number of points since  $g$  is a local homeomorphism and  $F_\ell$  is semi-algebraic. Let  $x_{\ell,1}, \dots, x_{\ell,p_\ell}$  denote these points. Note that if  $s_\ell^0 \subset Z$ , then  $s_{\ell'}^0 \subset Z$ ,  $s_\ell$  is a face of  $s_{\ell'}$  and  $x_{\ell,\lambda} \in F_\ell$ , then there exists a unique point  $x_{\ell',\lambda'} = \beta_{\ell,\ell'}(x_{\ell,\lambda}) \in F_{\ell'}$  such that  $\theta_\ell(s_\ell^0 \times \{x_{\ell,\lambda}\})$  is equal to the closure of  $(\theta_{\ell'}(s_{\ell'}^0 \times \{x_{\ell',\lambda'}\})) \cap g^{-1}(s_\ell^0)$ . Fix  $\ell$  and  $\lambda$  and set

$$V_{\ell,\lambda} = \bigcup \{ \theta_{\ell'}(s_{\ell'}^0 \times \{ \beta_{\ell,\ell'}(x_{\ell,\lambda}) \}) \mid s_{\ell'}^0 \subset Z \text{ and } s_\ell \text{ is a face of } s_{\ell'} \}.$$

By the previous remark,  $g \mid V_{\ell,\lambda}$  is a homeomorphism over the union of the  $s_{\ell'}^0 \subset Z$  such that  $s_\ell$  is a face of  $s_{\ell'}$ . The proposition is then proved, since the  $V_{\ell,\lambda}$  form a finite open cover of  $S$ .  $\square$

**Corollary 5.51.** *Let  $M$  be an  $\mathcal{S}^\infty$  submanifold of  $\mathbb{R}^k$  of dimension  $d$ . There exists a finite cover of  $M$  by semi-algebraic open sets  $M_i$  such that, for each  $M_i$ , one can find  $j_1, \dots, j_d \in \{1, \dots, k\}$  in such a way that the restriction to  $M_i$  of the projection mapping  $(x_1, \dots, x_k) \mapsto (x_{j_1}, \dots, x_{j_d})$  from  $\mathbb{R}^k$  onto  $\mathbb{R}^d$  is an  $\mathcal{S}^\infty$  diffeomorphism onto its image (stated differently, over each  $M_i$  one can express  $k - d$  coordinates as  $\mathcal{S}^\infty$  functions of the other  $d$  coordinates).*

**Proof:** Let  $\Pi: \mathbb{R}^k \rightarrow \mathbb{R}^d$  be the projection mapping that forgets the last  $k - d$  coordinates, and let  $M' \subset M$  be the set of points  $x$  such that  $\Pi$  induces an isomorphism from the tangent space  $T_x(M)$  onto  $\mathbb{R}^d$ . The function  $\Pi \mid M'$  is a local homeomorphism, hence, by Proposition 5.50, one can cover  $M'$  by the images of a finite number of semi-algebraic continuous sections (i.e. local inverses) of  $\Pi \mid M'$ , defined over semi-algebraic open sets of  $\mathbb{R}^d$ ; these sections are  $\mathcal{S}^\infty$  functions by Theorem 3.25 (Implicit Function Theorem). We do the same with projections onto all other  $k - d$ -coordinates, thereby exhausting the manifold.  $\square$

We now introduce a notion of local dimension.

**Proposition 5.52.** *Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set, and let  $a$  be a point of  $S$ . There exists a semi-algebraic open neighborhood  $U$  of  $x$  in  $S$  such that, for any other semi-algebraic open neighborhood  $U'$  of  $x$  in  $S$  contained in  $U$ , one has  $\dim(U) = \dim(U')$ .*

**Proof:** Clear by the properties of the dimension and Theorem 5.48 (local conic structure).  $\square$

Let  $S \subset \mathbb{R}^k$  be a semi-algebraic set and  $x$  a point of  $S$ . With  $U$  as in Proposition 5.52, one calls  $\dim(U)$  the **local dimension of  $S$  at  $x$** , denoted  $\dim(S_x)$ .

A point  $x \in S$  is a **smooth point of dimension  $d$**  of  $S$  if there exists a semi-algebraic open neighborhood  $U$  of  $\mathbb{R}^k$  such that  $S \cap U$  is an  $\mathcal{S}^\infty$  manifold of dimension  $d$ . Note that a smooth point of  $S$  of dimension  $d$  has local dimension  $d$ .

**Proposition 5.53.** *Let  $S$  be a semi-algebraic set of dimension  $d$ . There exists a non-empty semi-algebraic subset  $T \subset S$  such that every point of  $T$  is a smooth point of dimension  $d$  and  $S^{(d)} = \{x \in S \mid \dim(S_x) = d\}$  is a non-empty closed semi-algebraic subset  $S$ , which is the closure of  $T$ . Moreover  $\dim(S \setminus S^{(d)}) < d$ .*

**Proof:** By Theorem 5.38, the set  $S$  is a finite union of semi-algebraic sets  $S_i$ , each  $\mathcal{S}^\infty$  diffeomorphic to  $(0, 1)^{d(i)}$ . Let  $T$  be the union of the  $S_i$  such that  $d(i) = d$  (there are such  $S_i$  since  $d = \sup(d(i))$ ). It is clear that every point of  $T$  is a smooth point of dimension  $d$ . Let  $S'$  be the closure in  $S$  of  $T$ . Of course  $S' \subset S^{(d)}$ . If  $x \notin S'$ , there is a sufficiently small open neighborhood  $U$  of  $x$  such that  $S_i \cap U \neq \emptyset$  implies  $d(i) < d$  hence  $x \notin S^{(d)}$ . Therefore  $\mathbb{R}^k \setminus S^{(d)}$  is open. Note that  $S \setminus S^{(d)}$  contains no stratum of dimension  $d$ . This proves the claim.  $\square$

**Proposition 5.54.** *Let  $S$  be a semi-algebraic set. There exist non-empty semi-algebraic subsets of  $S$ ,  $S_1, \dots, S_\ell$  such that every point of  $S_i$  is a smooth point of dimension  $d(i)$  and  $S$  is the union of the closure of the  $S_i$ .*

**Proof:** The proof is by induction on the dimension of  $S$ . The claim is obviously true for  $\dim(S) = 1$ , since  $S$  is a finite number of points and intervals, and a closed interval is the closure of the corresponding open interval.

Suppose by induction hypothesis that the claim holds for all semi-algebraic sets of dimension  $< d$  and consider a semi-algebraic set  $S$  of dimension  $d$ . By Proposition 5.53, the set  $S^{(d)} = \{x \in S \mid \dim(S_x) = d\}$  is the closure of a semi-algebraic subset  $T_1$  such that every point of  $T_1$  is a smooth point of dimension  $d$ . Define  $S_1 = S \setminus S^{(d)}$ . It follows from Proposition 5.53 that the dimension of  $S_1$  is  $< d$ , and the claim follows by applying the induction hypothesis to  $S_1$ .  $\square$

## 5.9 Semi-algebraic Sard's Theorem

**Definition 5.55. [Critical point]** If  $f: N \rightarrow M$  is an  $\mathcal{S}^\infty$  function between two  $\mathcal{S}^\infty$  submanifolds  $N$  and  $M$ , then a **critical point** of  $f$  is a point  $x$  of  $N$  where the rank of the differential  $Df(x): T_x(N) \rightarrow T_{f(x)}(M)$  is strictly smaller than the dimension of  $M$ ; a **critical value** of  $f$  is the image of a critical point under  $f$ . A **regular point** of  $f$  on  $N$  is a point which is not critical and a **regular value** of  $f$  on  $M$  is a value of  $f$  which is not critical.  $\square$

We now give the semi-algebraic version of Sard's Theorem.

**Theorem 5.56. [Sard's theorem]** *Let  $f: N \rightarrow M$  be an  $\mathcal{S}^\infty$  function between two  $\mathcal{S}^\infty$  submanifolds. The set of critical values of  $f$  is a semi-algebraic subset of  $M$  whose dimension is strictly smaller than the dimension of  $M$ .*

The proof of Theorem 5.56 uses the constant rank theorem which can be proved from the inverse function theorem for  $\mathcal{S}^\infty$  functions.

**Theorem 5.57. [Constant Rank]** *Let  $f$  be a  $\mathcal{S}^\infty$  function from a semi-algebraic open set  $A$  of  $\mathbb{R}^k$  into  $\mathbb{R}^m$  such that the rank of the derivative  $df(x)$  is constant and equal to  $p$  over  $A$ , and  $a$  be a point in  $A$ .*

*There exists a semi-algebraic open neighborhood  $U$  of  $a$  which is contained in  $A$ , an  $\mathcal{S}^\infty$  diffeomorphism  $u: U \rightarrow (-1, 1)^k$ , a semi-algebraic open set  $V \supset f(U)$ , and an  $\mathcal{S}^\infty$  diffeomorphism  $v: (-1, 1)^m \rightarrow V$  such that  $f|_U = v \circ g \circ u$ , where  $g: (-1, 1)^k \rightarrow (-1, 1)^m$  is the mapping  $(x_1, \dots, x_k) \mapsto (x_1, \dots, x_p, 0, \dots, 0)$ .*

**Proof:** Without loss of generality let  $a = O$  be the origin and  $f(a) = O$ . Then,  $df(O): \mathbb{R}^k \rightarrow \mathbb{R}^m$  is a linear map of rank  $p$ . Let  $M \subset \mathbb{R}^k$  be the kernel of  $df(O)$  and  $N \subset \mathbb{R}^m$  be its image. It is clear that  $\dim(M) = k - p$  and  $\dim(N) = p$ .

Without loss of generality we can assume that  $M$  is spanned by the last  $k - p$  coordinates of  $\mathbb{R}^k$ , and we denote by  $M'$  the subspace spanned by the first  $p$  coordinates.

We also assume without loss of generality that  $N$  is spanned by the first  $p$  coordinates of  $\mathbb{R}^m$  and we denote by  $N'$  the subspace spanned by the last  $m - p$  coordinates. We will denote by  $\pi_M$  (resp.  $\pi_N, \pi_{N'}$ ) the projection maps from  $\mathbb{R}^k$  to  $M$  (resp.  $\mathbb{R}^m$  to  $N, N'$ ).

Let  $U'$  be a neighborhood of  $O$  in  $A$  and consider the map  $\tilde{u}_1: U' \rightarrow N \times M$  defined by

$$\tilde{u}_1(x) = (\pi_N(f(x)), \pi_M(x)).$$

Clearly,  $d\tilde{u}_1(O)$  is invertible, so by Proposition 3.24 (Inverse Function Theorem), there exists a neighborhood  $U''$  of the origin such that  $\tilde{u}_1|_{U''}$  is an  $\mathcal{S}^\infty$  diffeomorphism and  $\tilde{u}_1(U'')$  contains the set  $I_k(r) = (-r, r)^{k-p} \times (-r, r)^p$  for some sufficiently small  $r > 0$ . Let  $U = \tilde{u}_1^{-1}(I_k(r))$  and  $u_1 = \tilde{u}_1|_U$ .

Let  $V'$  be a neighborhood of the origin in  $\mathbb{R}^m$  containing  $f(U)$  and define  $\tilde{v}_1: V' \rightarrow N \times N'$  by

$$\tilde{v}_1(y) = (\pi_N(y), \pi_{N'}(y - f(u_1^{-1}(\pi_N(y), O)))).$$

Shrinking  $U$  and  $r$  if necessary, we can choose a neighborhood  $V'' \subset V'$  containing  $f(U)$ , such that  $\tilde{v}_1|_{V''}$  is an  $\mathcal{S}^\infty$  diffeomorphism. To see this observe that  $d\tilde{v}_1(O)$  is invertible, and apply Proposition 3.24 (Inverse Function Theorem). Shrink  $r$  if necessary so that  $\tilde{v}_1(V'')$  contains

$$I_m(r) = (-r, r)^p \times (-r, r)^{m-p}.$$

Let  $V = \tilde{v}_1^{-1}(I_m(r))$  and  $v_1 = \tilde{v}_1|_V$ . Finally, let  $u: U \rightarrow I_k(1)$  be defined by  $u(x) = u(x)/r$  and let  $v: I_m(1) \rightarrow V$  be the  $\mathcal{S}^\infty$  diffeomorphism defined by  $v(y) = v_1^{-1}(ry)$ .

We now prove that  $f|_U = v \circ g \circ u$ , where  $g: (-1, 1)^k \rightarrow (-1, 1)^m$  is the projection mapping  $(x_1, \dots, x_k) \mapsto (x_1, \dots, x_p, 0, \dots, 0)$ .

Since the rank of the derivative  $df(x)$  is constant and equal to  $p$  for all  $x \in U$ , we have that for each  $x \in U$  the image  $N_x$  of  $df(x)$  is a  $p$ -dimensional linear subspace of  $\mathbb{R}^m$ . Also, choosing  $r$  small enough we can ensure that  $\pi_N$  restricted to  $N_x$  is a bijection. We let  $L_x: N \rightarrow N_x$  denote the inverse of this bijection.

Now, consider the  $\mathcal{S}^\infty$  map  $f_1: (-r, r)^k \rightarrow \mathbb{R}^m$  defined by,

$$f_1(z_1, z_2) = f(u_1^{-1}(\pi_N(z_1), z_2)).$$

We first show that  $f_1(z_1, z_2)$  is in fact independent of  $z_2$ .

Clearly,

$$f(x) = f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))).$$

Differentiating using the chain rule, denoting  $g = du^{-1}(\pi_N(f(x)), \pi_M(x))$ , for all  $t \in \mathbb{R}^k$ ,

$$\begin{aligned} df(x)(t) &= d_1 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) \circ g \circ \pi_N \circ df(x)(t) \\ &\quad + d_2 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) \circ g \circ \pi_M(t), \end{aligned}$$

where  $d_i$  is the derivative with respect to  $z_i$ . Note also that,

$$df(x)(t) = L_x \circ \pi_N \circ df(x)(t).$$

Hence, with  $L = (L_x - d_1 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))))$

$$\begin{aligned} &d_2 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) \circ du^{-1}(\pi_N(f(x)), \pi_M(x)) \circ \pi_M(t) \\ &= L \circ du^{-1}(\pi_N(f(x)), \pi_M(x)) \circ \pi_N \circ df(x)(t). \end{aligned}$$

Let  $S_x$  denote the linear map

$$L_x - d_1 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) \circ du^{-1}(\pi_N(f(x)), \pi_M(x)): N \rightarrow N_x.$$

For  $t \in M'$ ,  $\pi_M(t) = 0$  and hence,  $S_x \circ \pi_N \circ df(x)(t) = 0$ . Since,  $\pi_N \circ df(x)$  is a bijection onto  $N$ , this implies that  $S_x = 0$ . Therefore, we get that

$$d_2 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) \circ du^{-1}(\pi_N(f(x)), \pi_M(x)) \circ \pi_M(t) = 0$$

for all  $t \in \mathbb{R}^k$  implying that

$$d_2 f_1(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) = 0$$

for all  $x \in U$ . This shows that  $f_1(z_1, z_2)$  is in fact independent of  $z_2$ .

Suppose now that  $v_1(y) \in N$  for some  $y \in V$ . This means that,

$$\pi_{N'}(y - f(u_1^{-1}(\pi_N(y), O))) = 0.$$

Let  $u_1^{-1}(\pi_N(y), O) = x$ . It follows from the definition of  $u_1$  and from our assumption that  $\pi_N(f(x)) = \pi_N(y)$  and  $\pi_{N'}(y) = \pi_{N'}(f(x))$ . Hence,  $y = f(x)$ .

Conversely, suppose that  $y = f(x)$ . Then using the fact that  $f_1(z_1, z_2)$  does not depend on  $z_2$  and the fact that  $u_1$  is injective we get that

$$f(u_1^{-1}(\pi_N(y), O)) = f(u_1^{-1}(\pi_N(f(x)), \pi_M(x))) = f(x) = y$$

and hence  $\pi_{N'}(y - f(u_1^{-1}(\pi_N(y), O))) = 0$ . Thus,  $v_1(y) \in N$ . □

**Proof of Theorem 5.56:** By Corollary 5.51, one may suppose that  $M$  is a semi-algebraic open set of  $\mathbb{R}^m$ . Let  $S \subset N$  be the set of critical points of  $f$ . The set  $S$  is semi-algebraic since the partial derivatives of  $f$  are  $\mathcal{S}^\infty$  functions. By Proposition 5.40,  $S$  is a finite union of semi-algebraic sets  $S_i$  that are the images of  $\mathcal{S}^\infty$  embeddings  $\varphi_i: (0, 1)^{d(i)} \rightarrow N$ . The rank of the composite function  $f \circ \varphi_i$  is  $< m$ . It remains to prove that the dimension of the image of  $f \circ \varphi_i$  is  $< m$ . This is done in the following lemma.

**Lemma 5.58.** *Let  $g: (0, 1)^d \rightarrow \mathbb{R}^m$  be an  $\mathcal{S}^\infty$  function such that the rank of the differential  $dg(x)$  is everywhere  $< m$ . Then, the dimension of the image of  $g$  is  $< m$ .*

**Proof of Lemma 5.58:** Let us suppose that  $\dim(g((0, 1)^d)) = m$ . By applying Corollary 5.51 to  $g$ , one can find a semi-algebraic open set  $U$  of  $\mathbb{R}^m$  that is contained in  $g((0, 1)^d)$  and a semi-algebraic homeomorphism  $\theta: U \times F \rightarrow g^{-1}(U)$  such that  $g \circ \theta$  is the projection of  $U \times F$  onto  $U$ . If  $x \in g^{-1}(U)$ , then the image under  $g$  of every semi-algebraic open neighborhood of  $x$  is a semi-algebraic open neighborhood of  $g(x)$  and is thus of dimension  $m$ . If for  $x$  one chooses a point where the rank of  $dg(x)$  is maximal (among the values taken over  $g^{-1}(U)$ ), then one obtains a contradiction with Theorem 5.57 (Constant Rank). □ □

## 5.10 Bibliographical Notes

The geometric technique underlying the cylindrical decomposition method can be found already in [160], for algebraic sets. The specific cylindrical decomposition method using subresultant coefficients comes from Collins [45].

Triangulation of semi-algebraic sets seems to appear for the first time in [28].

Hardt's triviality theorem appears originally in [83].

---

## Elements of Topology

In this chapter, we introduce basic concepts of algebraic topology adapted to semi-algebraic sets. We show how to associate to semi-algebraic sets discrete objects (the homology and cohomology groups) that are invariant under semi-algebraic homeomorphisms. In Section 6.1, we develop a combinatorial theory for homology and cohomology that applies only to simplicial complexes. In Section 6.2 we show how to extend this theory to closed semi-algebraic sets using the triangulation theorem proved in Chapter 5. In Section 6.3 we define homology groups, Borel-Moore homology groups and Euler-Poincaré characteristic for special cases of locally closed semi-algebraic sets.

### 6.1 Simplicial Homology Theory

#### 6.1.1 The Homology Groups of a Simplicial Complex

We define the simplicial homology groups of a simplicial complex  $K$  in a combinatorial manner. We use the notions and notation introduced in Section 5.6.

Given a simplicial complex  $K$ , let  $K_i$  be the set of  $i$ -dimensional simplices of  $K$ . In particular,  $K_0$  is the set of vertices of  $K$ .

##### 6.1.1.1 Chain Groups

Let  $p \in \mathbb{N}$ . A non-degenerate **oriented  $p$ -simplex** is a  $p$ -simplex  $[a_0, \dots, a_p]$  together with an equivalence class of total orderings on the set of vertices  $\{a_0, \dots, a_p\}$ , two orderings are equivalent if they differ by an even permutation of the vertices. Thus, a simplex has exactly two orientations. If  $a_0, \dots, a_p$  are not affinely independent, we set  $[a_0, \dots, a_p] = 0$ , which is a degenerate oriented  $p$ -simplex.



Abusing notation, if  $s = [a_0, \dots, a_p]$  is a  $p$ -simplex, we denote by  $[a_0, \dots, a_p]$  the oriented simplex corresponding to the order  $a_0 < a_1 < \dots < a_p$  on the vertices. So,  $s = [a_0, \dots, a_p]$  is an oriented simplex and  $-s = [a_1, a_0, a_2, \dots, a_p]$  is the oppositely oriented simplex.

Given a simplicial complex  $K$ , the  $\mathbb{Q}$ -vector space generated by the  $p$ -dimensional oriented simplices of  $K$  is called the  **$p$ -chain group** of  $K$  and is denoted  $C_p(K)$ . The elements of  $C_p(K)$  are called the  **$p$ -chains** of  $K$ . Notice that if  $K$  contains no  $p$ -simplices then  $C_p(K)$  is a  $\mathbb{Q}$ -vector space generated by the empty set, which is  $\{0\}$ . Since  $K_p$  is finite,  $C_p(K)$  is finite dimensional. An element of  $C_p(K)$  can be written  $c = \sum_i n_i s_i$ ,  $n_i \in \Lambda$ ,  $s_i \in K_p$ . For  $p < 0$ , we define  $C_p(K) = 0$ . When  $s$  is the oriented  $p$ -simplex  $[a_0, \dots, a_p]$ , we define  $[b, s]$  to be the oriented  $p+1$ -simplex  $[b, a_0, \dots, a_p]$ . If  $c = \sum_i n_i s_i$ , (with  $n_i \in \Lambda$ ) is a  $p$ -chain, then we define  $[b, c]$  to be  $\sum_i n_i [b, s_i]$ .

Given an oriented  $p$ -simplex  $s = [a_0, \dots, a_p]$ ,  $p > 0$ , the **boundary** of  $s$  is the  $(p-1)$ -chain

$$\partial_p(s) = \sum_{0 \leq i \leq p} (-1)^i [a_0, \dots, a_{i-1}, \hat{a}_i, a_{i+1}, \dots, a_p],$$

where the hat  $\hat{\phantom{a}}$  means that the corresponding vertex is omitted.

The map  $\partial_p$  extends linearly to a homomorphism  $\partial_p: C_p(K) \rightarrow C_{p-1}(K)$  by the rule

$$\partial_p\left(\sum_i n_i s_i\right) = \sum_i n_i \partial_p(s_i).$$

Note that, if  $c$  is a  $p$ -chain,  $\partial_{p+1}([b, c]) = c - [b, \partial_p(c)]$ .

For  $p \leq 0$ , we define  $\partial_p = 0$ . Thus, we have the following sequence of vector space homomorphisms,

$$\dots \longrightarrow C_p(K) \xrightarrow{\partial_p} C_{p-1}(K) \xrightarrow{\partial_{p-1}} C_{p-2}(K) \xrightarrow{\partial_{p-2}} \dots \xrightarrow{\partial_1} C_0(K) \xrightarrow{\partial_0} 0.$$

Using the definition of  $\partial_p$  and expanding, it is not too difficult to show that, for all  $p$

$$\partial_{p-1} \circ \partial_p = 0.$$

The sequence of pairs  $\{(C_p(K), \partial_p)\}_{p \in \mathbb{N}}$  is denoted  $\mathbf{C}_\bullet(K)$ .

Given two simplicial complexes  $K, L$ , a map  $\phi: |K| \rightarrow |L|$  is a **simplicial map** if it is the piecewise linear extension to each simplex of a map  $\phi_0: K_0 \rightarrow L_0$  that maps the vertices of every simplex in  $K$  to the vertices of a simplex in  $L$  (not necessarily of the same dimension). A simplicial map  $\phi$  defines a sequence of homomorphisms  $C_p(\phi)$  from  $C_p(K)$  to  $C_p(L)$  by

$$C_p(\phi)[a_0, \dots, a_p] = [\phi_0(a_0), \dots, \phi_0(a_p)].$$

Notice that the right hand side is automatically zero if  $\phi_0$  is not injective on the set  $\{a_0, \dots, a_p\}$ , in which case  $[\phi_0(a_0), \dots, \phi_0(a_p)]$  is a degenerate simplex. Also note that a simplicial map is automatically semi-algebraic.

### 6.1.1.2 Chain Complexes and Chain Maps

The chain groups obtained from a simplicial complex are a special case of more general abstract algebraic objects called chain complexes. The homomorphisms between the chain groups obtained from simplicial maps are then special cases of the more general chain homomorphisms, which we introduce below.

A sequence  $\{C_p\}$ ,  $p \in \mathbb{Z}$ , of vector spaces together with a sequence  $\{\partial_p\}$  of homomorphisms  $\partial_p: C_p \rightarrow C_{p-1}$  for which  $\partial_{p-1} \circ \partial_p = 0$  for all  $p$  is called a **chain complex**. Given two chain complexes,  $C_\bullet = (C_p, \partial_p)$  and  $C'_\bullet = (C'_p, \partial'_p)$ , a **chain homomorphism**  $\phi_\bullet: C_\bullet \rightarrow C'_\bullet$  is a sequence of homomorphisms  $\phi_p: C_p \rightarrow C'_p$  for which  $\partial'_p \circ \phi_p = \phi_{p-1} \circ \partial_p$  for all  $p$ .

In other words, the following diagram is commutative.

$$\begin{array}{ccccccc}
 \dots & \longrightarrow & C_p & \xrightarrow{\partial_p} & C_{p-1} & \longrightarrow & \dots \\
 & & \downarrow \phi_p & & \downarrow \phi_{p-1} & & \\
 \dots & \longrightarrow & C'_p & \xrightarrow{\partial'_p} & C'_{p-1} & \longrightarrow & \dots
 \end{array}$$

Notice that if  $\phi: K \rightarrow K'$  is a simplicial map, then  $C_\bullet(\phi): C_\bullet(K) \rightarrow C_\bullet(K')$  is a chain homomorphism between the chain complexes  $C_\bullet(K)$  and  $C_\bullet(K')$ .

### 6.1.1.3 Homology of Chain Complexes

Given a chain complex  $C_\bullet$ , the elements of  $B_p(C_\bullet) = \text{Im}(\partial_{p+1})$  are called  **$p$ -boundaries** and those of  $Z_p(C_\bullet) = \text{Ker}(\partial_p)$  are called  **$p$ -cycles**. Note that, since  $\partial_{p-1} \circ \partial_p = 0$ ,  $B_p(C_\bullet) \subset Z_p(C_\bullet)$ . The **homology groups**  $H_p(C_\bullet)$  are defined by  $H_p(C_\bullet) = Z_p(C_\bullet) / B_p(C_\bullet)$ .

Note that, by our definition, the homology groups  $H_p(C_\bullet)$  are all  $\mathbb{Q}$ -vector spaces (finite dimensional if the vector spaces  $C_p$ 's are themselves finite dimensional). We still refer to them as groups as this is standard terminology in algebraic topology where more general rings of coefficients, for instance the integers, are often used in the definition of the chain complexes. In such situations, the homology groups are not necessarily vector spaces over a field, but rather modules over the corresponding ring.

This sequence of groups together with the sequence of homomorphisms which sends each  $H_p(C_\bullet)$  to  $0 \in H_{p-1}(C_\bullet)$  constitutes a chain complex  $(H_p(C_\bullet), 0)$  which is denoted by  $H_\star(C_\bullet)$ .

**Lemma 6.1.** *Given two chain complexes  $C_\bullet$  and  $C'_\bullet$ , a chain homomorphism  $\phi_\bullet: C_\bullet \rightarrow C'_\bullet$  induces a homomorphism  $H_\star(\phi_\bullet): H_\star(C_\bullet) \rightarrow H_\star(C'_\bullet)$  which respects composition. In other words, given another chain homomorphism  $\psi_\bullet: C'_\bullet \rightarrow C''_\bullet$ ,*

$$H_\star(\psi_\bullet \circ \phi_\bullet) = H_\star(\psi_\bullet) \circ H_\star(\phi_\bullet)$$

and  $H_\star(\text{Id}_{C_\bullet}) = \text{Id}_{H_\star(C_\bullet)}$ .

**Proof:** Using the fact that the diagram of a chain homomorphism commutes, we see that a chain homomorphism carries cycles to cycles and boundaries to boundaries. Thus, the chain homomorphism  $\phi$  induces homomorphisms

$$Z_p(\phi_\bullet): Z_p(C_\bullet) \rightarrow Z_p(C'_\bullet),$$

$$B_p(\phi_\bullet): B_p(C_\bullet) \rightarrow B_p(C'_\bullet).$$

Thus, it also induces a homomorphism

$$H_p(\phi_\bullet): H_p(C_\bullet) \rightarrow H_p(C'_\bullet).$$

The remaining claims follow easily.  $\square$

#### 6.1.1.4 Homology of a Simplicial Complex

**Definition 6.2.** Given a simplicial complex  $K$ ,  $H_p(K) = H_p(C_\bullet(K))$  is the  $p$ -th simplicial homology group of  $K$ . As a special case of Lemma 6.1, it follows that a simplicial map from  $K$  to  $L$  induces homomorphisms between the homology groups  $H_p(K)$  and  $H_p(L)$ .

We denote by  $H_\star(K)$  the chain complex  $(H_p(K), 0)$  and call it the **homology of  $K$** .

It is clear from the definition that  $H_p(K)$  is a finite dimensional  $\mathbb{Q}$ -vector space. The dimension of  $H_p(K)$  as a  $\mathbb{Q}$ -vector space is called the  $p$ -th **Betti number** of  $K$  and denoted  $b_p(K)$ .

The **Euler-Poincaré characteristic** of  $K$  is

$$\chi(K) = \sum_i (-1)^i b_i(K). \quad \square$$

**Proposition 6.3.** *Let  $n_i(K)$  be the number of simplexes of dimension  $i$  of  $K$ . Then*

$$\chi(K) = \sum_i (-1)^i n_i(K).$$

**Proof:** Recall from the definition of  $H_i(K)$  that,

$$b_i(K) = \dim H_i(K) = \dim \text{Ker}(\partial_i) - \dim \text{Im}(\partial_{i+1}).$$

Moreover,

$$n_i(K) = \dim C_i(K) = \dim \text{Ker}(\partial_i) + \dim \text{Im}(\partial_i).$$

An easy calculation now shows that,

$$\begin{aligned} \chi(K) &= \sum_i (-1)^i b_i(K) \\ &= \sum_i (-1)^i (\dim \text{Ker}(\partial_i) - \dim \text{Im}(\partial_{i+1})) \\ &= \sum_i (-1)^i (\dim \text{Ker}(\partial_i) + \dim \text{Im}(\partial_i)) \\ &= \sum_i (-1)^i n_i. \end{aligned}$$

□

### 6.1.2 Simplicial Cohomology Theory

We have defined the homology groups of a simplicial complex  $K$  in the previous section. We now define a dual notion – namely that of cohomology groups. One reason for defining cohomology groups is that in many situations, it is more convenient and intuitive to reason with the cohomology groups than with the homology groups.

Given a simplicial complex  $K$ , we will denote by  $C^p(K)$  the vector space dual to  $C_p(K)$ , and the sequence of homomorphisms,

$$0 \rightarrow C^0(K) \xrightarrow{\delta^0} C^1(K) \xrightarrow{\delta^1} C^2(K) \cdots C^p(K) \xrightarrow{\delta^p} C^{p+1}(K) \xrightarrow{\delta^{p+1}} \cdots$$

is called the **cochain complex** of  $K$ . Here,  $\delta^p$  is the homomorphism dual to  $\partial_{p+1}$  in the chain complex  $C^\bullet(K)$ . The sequence of pairs  $\{(C^p(K), \delta^p)\}_{p \in \mathbb{N}}$  is denoted by  $C^\bullet(K)$ . Notice that each  $\phi \in C^p(K)$  is a linear functional on the vector space  $C^p(K)$ , and thus  $\phi$  is determined by the values it takes on each  $i$ -simplex of  $K$ .

#### 6.1.2.1 Cochain Complexes

The dual notion for chain complexes is that of cochain complexes. A sequence  $\{C^p\}$ ,  $p \in \mathbb{Z}$ , of vector spaces together with a sequence  $\{\delta^p\}$  of homomorphisms  $\delta^p: C^p \rightarrow C^{p+1}$  for which  $\delta^{p+1} \circ \delta^p = 0$  for all  $p$  is called a **cochain complex**.

Given two cochain complexes,  $C^\bullet = (C^p, \delta^p)$  and  $D^\bullet = (D^p, \delta'^p)$ , a **cochain homomorphism**  $\phi^\bullet: C^\bullet \rightarrow D^\bullet$  is a sequence of homomorphisms  $\phi^p: C^p \rightarrow D^p$  for which  $\delta'^p \circ \phi^p = \phi^{p+1} \circ \delta^p$  for all  $p$ .

In other words, the following diagram is commutative.

$$\begin{array}{ccccccc}
 \dots & \longrightarrow & C^p & \xrightarrow{\delta^p} & C^{p+1} & \longrightarrow & \dots \\
 & & \downarrow \phi^p & & \downarrow \phi^{p+1} & & \\
 \dots & \longrightarrow & D^p & \xrightarrow{\delta'^p} & D^{p+1} & \longrightarrow & \dots
 \end{array}$$

It is clear that given a chain complex  $C_\bullet = \{(C_p, \partial_p)\}_{p \in \mathbb{N}}$ , we can obtain a corresponding cochain complex  $C^\bullet = \{(C^p, \delta^p)\}_{p \in \mathbb{N}}$  by taking duals of each term and homomorphisms. Doing so reverses the direction of every arrow in the corresponding diagram.

### 6.1.2.2 Cohomology of Cochain Complexes

The elements of  $B^p(C^\bullet) = \text{Im}(\delta^{p-1})$  are called  $p$ -**coboundaries** and those of  $Z^p(C^\bullet) = \text{Ker}(\delta^p)$  are called  $p$ -**cocycles**. It is easy to verify that  $B^p(C^\bullet) \subset Z^p(C^\bullet)$ . The **cohomology groups**,  $H^p(C^\bullet)$ , are defined by

$$H^p(C^\bullet) = \frac{Z^p(C^\bullet)}{B^p(C^\bullet)}.$$

This sequence of groups together with the sequence of homomorphisms which sends each  $H^p(C^\bullet)$  to  $0 \in H^{p+1}(C^\bullet)$  constitutes a chain complex  $(H^p(C^\bullet), 0)$  which is denoted by  $H^*(C^\bullet)$ .

It is an easy exercise in linear algebra to check that:

**Proposition 6.4.** *Let  $C_\bullet$  be a chain complex and  $C^\bullet$  the corresponding cochain complex. Then, for every  $p \geq 0$ ,  $H_p(C_\bullet) \cong H^p(C^\bullet)$ .*

Given a simplicial complex  $K$ , the  $p$ -th cohomology group  $H^p(C^\bullet(K))$  will be denoted by  $H^p(K)$ . The cohomology group  $H^0(K)$  has a particularly natural interpretation. It is the vector space of locally constant functions on  $|K|$ .

**Proposition 6.5.** *Let  $K$  be a simplicial complex such that  $K_0 \neq \emptyset$ . The cohomology group  $H^0(K)$  is the vector space of locally constant functions on  $|K|$ . As a consequence, the number of connected components of  $|K|$  is  $b_0(K)$ .*

**Proof:** Clearly,  $H^0(K)$  depends only on the 1-skeleton of  $K$ , that is the subcomplex of  $K$  consisting of the zero and one-dimensional simplices.

Let  $z \in C^0(K)$  be a cocycle, that is such that  $d^0(z) = 0$ . This implies that for any  $e = [u, v] \in K_1$ ,  $z(u) - z(v) = 0$ . Hence  $z$  takes a constant value on vertices in a connected component of  $|K|$ . Since  $B^0(C^\bullet(K))$  is 0, this shows that  $H^0(K)$  is the vector space of locally constant functions on  $|K|$ .

Using Proposition 6.4, the last part of the claim follows since the dimension of the vector space of locally constant functions on  $|K|$  is the number of connected components of  $|K|$ . □

It follows immediately from Lemma 6.1 that

**Lemma 6.6.** *Given two cochain complexes  $C^\bullet$  and  $C'^\bullet$ , a cochain homomorphism  $\phi^\bullet: C^\bullet \rightarrow C'^\bullet$  induces a homomorphism  $H^*(\phi^\bullet): H^*(C^\bullet) \rightarrow H^*(C'^\bullet)$  which respects composition. In other words, given another chain homomorphism  $\psi^\bullet: C'^\bullet \rightarrow C''^\bullet$ ,*

$$H^*(\psi^\bullet \circ \phi^\bullet) = H^*(\psi^\bullet) \circ H^*(\phi^\bullet) \text{ and } H^*(\text{Id}_{C^\bullet}) = \text{Id}_{H^*(C^\bullet)}.$$

**6.1.3 A Characterization of  $H^1$  in a Special Case.**

Let  $A$  be a simplicial complex and  $A^1, \dots, A^s$  sub-complexes of  $A$  such that, each  $A^i$  is connected and

$$A = A^1 \cup \dots \cup A^s, \\ H^1(A^i) = 0, 1 \leq i \leq s.$$

For  $1 \leq i < j < \ell \leq s$ , we denote by  $A^{ij}$  the sub-complex  $A^i \cap A^j$ , and by  $A^{ij\ell}$  the sub-complex  $A^i \cap A^j \cap A^\ell$ . We will denote by  $C_\alpha^{ij}$  the sub-complexes corresponding to connected components of  $A^{ij}$ , and by  $C_\beta^{ij\ell}$  the sub-complexes corresponding to connected components of  $A^{ij\ell}$ .

We will show that the simplicial cohomology group,  $H^1(A)$ , is isomorphic to the first cohomology group of a certain complex defined in terms of  $H^0(A^i)$ ,  $H^0(A^{ij})$  and  $H^0(A^{ij\ell})$ . This result will be the basis of an efficient algorithm for computing the first Betti number of semi-algebraic sets, which will be developed in Chapter 16.

Let

$$N^\bullet = N^0 \longrightarrow N^1 \longrightarrow N^2 \rightarrow 0$$

denote the complex

$$C^0(A) \xrightarrow{d^0} C^1(A) \xrightarrow{d^1} C^2(A) \rightarrow 0.$$

Note that  $N^\bullet$  is just a truncated version of the cochain complex of  $A$ . The coboundary homomorphisms  $d^0, d^1$  are identical to the ones in  $C^\bullet(A)$ .

For each  $h \geq 0$ , we define

$$\delta_0: \bigoplus_{1 \leq i \leq s} C^h(A^i) \longrightarrow \bigoplus_{1 \leq i < j \leq s} C^h(A^{ij})$$

as follows.

For  $\phi \in \bigoplus_{1 \leq i \leq s} C^0(A^i)$ ,  $1 \leq i < j \leq s$ , and each oriented  $h$ -simplex  $\sigma \in A_h^{ij}$ ,

$$\delta_0^h \phi_{i,j}(\sigma) = \phi_i(\sigma) - \phi_j(\sigma).$$

Similarly, we define

$$\delta_1^h: \bigoplus_{1 \leq i < j \leq s}^h C^h(A^{ij}) \longrightarrow \bigoplus_{1 \leq i < j < \ell \leq s} C^h(A^{ij\ell})$$

by defining for  $\psi \in \bigoplus_{1 \leq i < j \leq s} C^h(A^{ij})$ ,  $1 \leq i < j < \ell \leq s$  and each oriented  $h$ -simplex  $\sigma \in A_h^{ij}$ ,

$$(\delta_1^h \psi)_{ij\ell}(\sigma) = \psi_{j\ell}(\sigma) - \psi_{i\ell}(\sigma) + \psi_{ij}(\sigma).$$

Let

$$M^\bullet = M^0 \longrightarrow M^1 \longrightarrow M^2 \rightarrow 0$$

denote the complex

$$\bigoplus_{1 \leq i \leq s} C^0(A^i) \xrightarrow{D_0} \bigoplus_{1 \leq i \leq s} C^1(A^i) \oplus \bigoplus_{1 \leq i < j \leq s} C^0(A^{ij}) \xrightarrow{D_1} \bigoplus_{\ell+n=2} \bigoplus_{J_n} C^\ell(A^{i_1 \dots i_n}) \rightarrow 0.$$

where  $J_n = \{i_1 \dots i_n \mid 1 \leq i_1 < \dots < i_n \leq s\}$

The homomorphism  $D_0$  is defined by

$$D_0(\phi) = d^0(\phi) \oplus \delta_0^0(\phi), \phi \in \bigoplus_{1 \leq i \leq s} C^0(A^i)$$

and  $D_1$  is defined by

$$D_1(\phi \oplus \psi) = d^1(\phi) \oplus (-\delta_0^1(\phi) + d^0(\psi)) \oplus -\delta_1^0(\psi),$$

$$\phi \in \bigoplus_{1 \leq i \leq s} C^1(A^i), \psi \in \bigoplus_{1 \leq i < j \leq s} C^0(A^{ij}).$$

Finally let

$$L^\bullet = L^0 \longrightarrow L^1 \longrightarrow L^2 \rightarrow 0$$

denote the complex

$$\bigoplus_{1 \leq i \leq s} H^0(A^i) \xrightarrow{\delta_0} \bigoplus_{1 \leq i < j \leq s} H^0(A^{ij}) \xrightarrow{\delta_1} \bigoplus_{1 \leq i < j < \ell \leq s} H^0(A^{ij\ell}) \rightarrow 0.$$

Recall that  $H^0(X)$  can be identified as the vector space of locally constant functions on the simplicial complex  $X$ , and is thus a vector space whose dimension equals the number of connected components of  $X$ . The homomorphisms  $\delta_i$  in the complex  $L^\bullet$  are generalized restriction homomorphisms. Thus, for  $\phi \in \bigoplus_{1 \leq i \leq s} H^0(A^i)$ ,

$$(\delta_0 \phi)_{ij} = \phi_i|_{A^{ij}} - \phi_j|_{A^{ij}}$$

and for  $\psi \in \bigoplus_{1 \leq i < j \leq s} H^0(A^{ij})$ ,

$$(\delta_1 \psi)_{ij\ell} = \psi_{j\ell}|_{A^{ij\ell}} - \psi_{i\ell}|_{A^{ij\ell}} + \psi_{ij}|_{A^{ij\ell}}.$$

We now define a homomorphism of complexes,

$$F^\bullet: L^\bullet \rightarrow M^\bullet,$$

as follows:

For  $\phi \in L^0$  and  $u \in A_0^i$ ,

$$F^1(\phi)_i(u) = \phi_i(A^i).$$

For  $\psi \in L^2$  and  $e \in A_1^{ij}$ ,

$$F^2(\psi)_i = 0,$$

and

$$F^2(\psi)_{ij}(e) = \psi_{ij}(C_\alpha^{ij}),$$

where  $C_\alpha^{ij}$  is the connected component of  $A^{ij}$  containing  $e$ .

For  $\theta \in L^3$  and  $\sigma \in A_2^{ij\ell}$

$$F^3(\theta)_i = F^3(\theta)_{ij} = 0,$$

and

$$F^3(\theta)_{ij\ell} = \psi_{ij\ell}(C_{ij\ell}^\beta),$$

where  $C_{ij\ell}^\beta$  is the connected component of  $A^{ij\ell}$  containing  $\sigma$ . It is easy to verify that  $F^\bullet$  is a homomorphism of complexes, and thus induces an homomorphism

$$H^*(F^\bullet): H^*(L^\bullet) \rightarrow H^*(M^\bullet).$$

We now prove that,

**Proposition 6.7.** *The induced homomorphism,*

$$H^1(F^\bullet): H^1(L^\bullet) \rightarrow H^1(M^\bullet)$$

*is an isomorphism.*

**Proof:** We first prove that,  $H^1(F^\bullet): H^1(L^\bullet) \rightarrow H^1(M^\bullet)$  is surjective. Let  $z = \phi \oplus \psi \in M^1$  be a cocycle, where

$$\phi \in \bigoplus_{1 \leq i \leq s} C^1(A^i)$$

and

$$\psi \in \bigoplus_{1 \leq i < j \leq s} C^0(A^{ij})$$

Since  $z$  is a cocycle, that is  $D_0(z) = 0$ , we have from the definition of  $D_0$  that,

$$\begin{aligned} d^1 \phi &= 0, \\ \delta_1^0 \psi &= 0, \\ \delta_0^1 \phi + d^0 \psi &= 0. \end{aligned}$$



From the first property, and the fact  $H^1(A^i) = 0$  for each  $i, 1 \leq i \leq \ell$ , we deduce that there exists

$$\theta \in M^0 = \bigoplus_{1 \leq i \leq \ell} C^0(A^i)$$

such that, for  $e = [u, v] \in A_1^i$ ,

$$\phi_i(e) = \theta_i(u) - \theta_i(v). \quad (6.1)$$

As a consequence of the second property, we have that for  $1 \leq i < j < \ell \leq s$ , and  $u \in A_0^{ij\ell}$ ,

$$\psi_{j\ell}(u) - \psi_{i\ell}(u) + \psi_{ij}(u) = 0. \quad (6.2)$$

Finally, from the third property we get that, for  $1 \leq i < j \leq s$  and  $e = [u, v] \in A_1^{ij}$ , and  $\theta$  defined above,

$$\theta_i(u) - \theta_i(v) - \theta_j(u) + \theta_j(v) + \psi_{ij}(u) - \psi_{ij}(v) = 0. \quad (6.3)$$

We now define  $z' = 0 \oplus \gamma \in M^1$  by defining, for  $1 \leq i < j \leq s$  and  $u \in A_0^{ij}$ ,

$$\gamma_{ij}(u) = \theta_i(u) - \theta_j(u) + \psi_{ij}(u).$$

From (6.3) it follows that  $\gamma_{ij}(u)$  is constant for all vertices  $u$  in any connected component of  $A^{ij}$ . Thus,  $z' \in F^1(L^1)$ . Next, for  $1 \leq i < j \leq s$ , and  $e = [u, v] \in A_1^{ij}$

$$\begin{aligned} d\gamma_{ij}(e) &= \theta_i(u) - \theta_j(u) + \psi_{ij}(u) - (\theta_i(v) - \theta_j(v) + \psi_{ij}(v)) \\ &= 0. \end{aligned}$$

where we again use (6.3). This shows that  $z'$  is a cycle.

Finally, it is easy to check, using the facts that  $\psi - \gamma = \delta\theta$ , and  $\phi = d^0\theta$ , that,  $z - z' = (d^0 + \delta_0^0)\theta$  is a coboundary in  $M^\bullet$ . This proves the surjectivity of  $H^1(F^\bullet)$ .

We now prove that  $H^1(F^\bullet)$  is injective by proving that for any  $z \in L^1$ , if  $F^1(z)$  is a coboundary in  $M^1$  then  $z$  must be a coboundary in  $L^1$ . Let

$$F^1(z) = (d^0 + \delta_0^0)\theta$$

for  $\theta \in \bigoplus_{1 \leq i \leq s} C^0(A^i)$ . We define  $\gamma \in \bigoplus_{1 \leq i \leq s} H^0(A^i)$  by defining

$$\gamma_i(A_i) = \theta_i(u)$$

for some  $u \in A_0^i$ . This is well defined, since by assumption each  $A^i$  is connected, and  $d^0\theta = 0$ , and thus we have that for each  $e = [u, v] \in A_1^i$ ,  $\theta_i(u) - \theta_i(v) = 0$ .

For any  $1 \leq i < j \leq s$ ,  $C_\alpha^{ij}$  a connected component of  $A^{ij}$ , and  $u \in C_{\alpha,0}^{ij}$ , we have that,

$$\begin{aligned} F^1(z)_{ij}(u) &= z_{ij}(C_\alpha^{ij}) \\ &= \theta_i(u) - \theta_j(u). \end{aligned}$$

It is now easy to check that  $\delta_0\gamma = z$ , proving that  $z$  is a coboundary in  $L^1$ . This proves the injectivity of  $H^1(F)$ . □

We now define a homomorphism of complexes,  $G^\bullet: N^\bullet \rightarrow M^\bullet$ , as follows.

First observe that for  $1 \leq i < j < \ell \leq s$ , there are natural restriction homomorphisms,

$$r_i^\bullet: C^\bullet(A) \rightarrow C^\bullet(A^i),$$

For  $\phi \in C^0(A)$ ,

$$G^0(\phi) = \bigoplus_{1 \leq i \leq s} \gamma_i^0(\phi_i).$$

For  $\psi \in C^1(A)$ ,

$$G^1(\psi) = \bigoplus_{1 \leq i \leq s} \gamma_i^1(\psi).$$

For  $\nu \in C^2(A)$ ,

$$G^2(\nu) = \bigoplus_{1 \leq i \leq s} \gamma_i^2(\nu).$$

We now prove that,

**Proposition 6.8.** *The induced homomorphism,*

$$H^1(G^\bullet): H^1(N^\bullet) \rightarrow H^1(M^\bullet)$$

*is an isomorphism.*

**Proof:** We first prove that  $H^1(G^\bullet)$  is surjective.

Let  $z = \phi \oplus \psi \in M^1$  be a cocycle, where

$$\phi \in \bigoplus_{1 \leq i \leq s} C^1(A^i), \quad \psi \in \bigoplus_{1 \leq i < j \leq s} C^0(A^{ij}).$$

Since  $z$  is a cocycle, that is  $D_0(z) = 0$ , we have from the definition of  $D_0$  that,

$$\begin{aligned} d^1\phi &= 0, \\ \delta_1^0\psi &= 0, \\ \delta_0^1\phi + d^0\psi &= 0. \end{aligned}$$

For  $1 \leq i < j < \ell \leq s$ , and  $u \in A_0^{ij\ell}$ ,

$$\psi_{j\ell}(u) - \psi_{i\ell}(u) + \psi_{ij}(u) = 0. \tag{6.4}$$

We now define  $\theta \in \bigoplus_{1 \leq i \leq s} C^0(A^i)$  such that,  $\delta_0(\theta) = \psi$ .

For  $1 \leq i \leq s$  and  $u \in A_0^i$  we define,

$$\theta_i(u) = \frac{1}{n_u} \sum_{\substack{1 \leq j \leq s \\ j \neq i, u \in A_0^j}} (-1)^{i-j} \psi_{ij}(u),$$

where  $n_u = \#\{j \mid u \in A_0^j\}$ .

It is easy to check using (6.4) that  $\delta_0(\theta) = \psi$ . Now define,  $z' = (\phi - d^0\theta) \oplus 0$ . Now,  $z' \in G^1(N^1)$ , since for  $1 \leq i < j \leq s$ , and  $e = [u, v] \in A_1^{ij}$ ,

$$\begin{aligned} (\phi - d^0\theta)_i(e) - (\phi - d^0\theta)_j(e) &= \phi_i(e) - \phi_j(e) - (\theta_i - \theta_j)(u - v) \\ &= (\phi_i(e) - \phi_j(e)) - \psi_{ij}(u) + \psi_{ij}(v) \\ &= (\delta_0^1\phi - d^0\psi)_{ij}(e) \\ &= 0. \end{aligned}$$

Also,  $z - z' = d^0\theta \oplus \psi = (d^0 + \delta_1^0)\theta$  is a coboundary. This show that  $H^1(G^\bullet)$  is surjective.

Finally, since  $G^1$  is obviously injective, it is clear that if the image of  $z \in N^1$  is a coboundary in  $M^1$ , then it must also be a coboundary in  $N^1$ , which shows that  $H^1(G^\bullet)$  is injective as well.  $\square$

We are now in a position to prove,

**Theorem 6.9.** *Let  $A$  be a simplicial complex and  $A^1, \dots, A^k$  sub-complexes of  $A$  such that, each  $A_i$  is connected and*

$$\begin{aligned} A &= A^1 \cup \dots \cup A^s, \\ H^1(A^i) &= 0, 1 \leq i \leq s, \end{aligned}$$

and let  $L^\bullet$  be the complex defined above. Then,

$$H^1(A) \cong H^1(L^\bullet).$$

**Proof:** The theorem follows directly from Proposition 6.7 and Proposition 6.8 proved above.  $\square$

### 6.1.4 The Mayer-Vietoris Theorem

In the next chapter, we will use heavily certain relations between the homology groups of two semi-algebraic sets and those of their unions and intersections. We start by indicating similar relations between the homology groups of the unions and intersections of sub-complexes of a simplicial complex. It turns out to be convenient to formulate these relations in terms of exact sequences.

A sequence of vector space homomorphisms,

$$\dots \xrightarrow{\phi_{i-2}} F_{i-1} \xrightarrow{\phi_{i-1}} F_i \xrightarrow{\phi_i} F_{i+1} \xrightarrow{\phi_{i+1}} \dots$$

is **exact** if and only if  $\text{Im}(\phi_i) = \text{Ker}(\phi_{i+1})$  for each  $i$ .

Let  $C_\bullet, C'_\bullet, C''_\bullet$  be chain complexes, and let  $\phi_\bullet: C_\bullet \rightarrow C'_\bullet, \psi_\bullet: C'_\bullet \rightarrow C''_\bullet$  be chain homomorphisms. We say that the sequence  $0 \rightarrow C_\bullet \xrightarrow{\phi_\bullet} C'_\bullet \xrightarrow{\psi_\bullet} C''_\bullet \rightarrow 0$  is a **short exact sequence of chain complexes** if in each dimension  $p$  the sequence  $0 \rightarrow C_p \xrightarrow{\phi_p} C'_p \xrightarrow{\psi_p} C''_p \rightarrow 0$  is an exact sequence of vector spaces.

We need the following lemma from homological algebra.

**Lemma 6.10. [Zigzag Lemma]** *Let  $0 \rightarrow C_\bullet \xrightarrow{\phi} C'_\bullet \xrightarrow{\psi} C''_\bullet \rightarrow 0$  be a short exact sequence of chain complexes. Then, there exist connecting homomorphisms,  $H_p(\partial)$ , making the following sequence exact*

$$\dots \xrightarrow{H_p(\phi_\bullet)} H_p(C'_\bullet) \xrightarrow{H_p(\psi_\bullet)} H_p(C''_\bullet) \xrightarrow{H_p(\partial)} H_{p-1}(C_\bullet) \xrightarrow{H_{p-1}(\phi_\bullet)} H_{p-1}(C'_\bullet) \dots$$

**Proof:** The proof is by “chasing” the following diagram:

$$\begin{array}{ccccccccc} 0 & \rightarrow & C_{p+1} & \xrightarrow{\phi_{p+1}} & C'_{p+1} & \xrightarrow{\psi_{p+1}} & C''_{p+1} & \rightarrow & 0 \\ & & \downarrow \partial_{p+1} & & \downarrow \partial'_{p+1} & & \downarrow \partial''_{p+1} & & \\ 0 & \rightarrow & C_p & \xrightarrow{\phi_p} & C'_p & \xrightarrow{\psi_p} & C''_p & \rightarrow & 0 \\ & & \downarrow \partial_p & & \downarrow \partial'_p & & \downarrow \partial''_p & & \\ 0 & \rightarrow & C_{p-1} & \xrightarrow{\phi_{p-1}} & C'_{p-1} & \xrightarrow{\psi_{p-1}} & C''_{p-1} & \rightarrow & 0. \end{array}$$

We have already seen in the proof of Lemma 6.1 that the chain homomorphisms  $\phi_\bullet$  and  $\psi_\bullet$  actually take boundaries to boundaries and hence the homomorphisms  $\phi_p$  (resp.  $\psi_p$ ) descend to homomorphisms on the homology vector spaces  $H_p(C_\bullet) \rightarrow H_p(C'_\bullet)$  (resp.,  $H_p(C'_\bullet) \rightarrow H_p(C''_\bullet)$ ). We denote these homomorphisms by  $H_p(\phi_\bullet)$  (resp.  $H_p(\psi_\bullet)$ ).

We now define the homomorphism  $H_p(\partial): H_p(C''_\bullet) \rightarrow H_{p-1}(C_\bullet)$ . Let  $\alpha''$  be a cycle in  $C''_p$ . By the exactness of the second row of the diagram we know that  $\psi_p$  is surjective and thus there exists  $\alpha' \in C'_p$  such that  $\psi_p(\alpha') = \alpha''$ . Let  $\beta' = \partial'_p(\alpha')$ .

We show that  $\beta' \in \text{Ker}(\psi_{p-1})$ . Using the commutativity of the diagram, we have that

$$\psi_{p-1}(\beta') = \psi_{p-1}(\partial'_p(\alpha')) = \partial''_p(\psi_p(\alpha')) = \partial''_p(\alpha'') = 0,$$

the last equality by virtue of the fact that  $\alpha''$  is a cycle.

By the exactness of the third row of the diagram, we have that  $\text{Im}(\phi_{p-1}) = \text{Ker}(\psi_{p-1})$  and hence  $\phi_{p-1}$  is injective. Thus, there exists a unique  $\beta \in C_{p-1}$  such that  $\beta' = \phi_{p-1}(\beta)$ . Moreover,  $\beta$  is a cycle in  $C_{p-1}$ . To see this, observe that

$$\phi_{p-2}(\partial_{p-1}(\beta)) = \partial'_{p-1}(\phi_{p-1}(\beta)) = \partial'_{p-1}(\beta') = \partial'_{p-1}(\partial'_p(\beta')) = 0.$$

Since  $\phi_{p-2}$  is injective, it follows that  $\partial_{p-1}(\beta) = 0$ , whence  $\beta$  is a cycle. Define  $H_p(\partial)(\overline{\alpha''}) = \overline{\beta}$ , where  $\overline{\beta}$  represents the homology class of the cycle  $\beta$ .

We now check that  $H_p(\partial)$  is a well-defined homomorphism and that the long sequence of homology is exact as claimed.

We first prove that the map defined above indeed is a well-defined homomorphism  $H_p(C''_\bullet) \rightarrow H_{p-1}(C_\bullet)$ . We first check that the homology class  $\bar{\beta}$  does not depend on the choice of  $\alpha' \in C'_p$  used in its definition. Let  $\alpha'_1 \in C'_p$  be such that  $\psi_p(\alpha'_1) = \alpha''$ . Let  $\beta'_1 = \partial'_p(\alpha'_1)$ . Now,  $\beta'_1$  is also in  $\text{Ker}(\psi_{p-1})$  and by the exactness of the third row of the diagram, there exists a unique cycle  $\beta_1 \in C_{p-1}$  such that  $\beta'_1 = \phi_{p-1}(\beta_1)$ .

Now,  $\alpha' - \alpha'_1 \in \text{Ker}(\psi_p)$ . Hence, there exists  $\alpha_0 \in C_p$  such that  $\phi_p(\alpha_0) = \alpha' - \alpha'_1$ , and using the commutativity of the diagram and the fact that  $\phi_{p-1}$  is injective, we have that  $\partial_p(\alpha_0) = \beta - \beta'_1$ , whence  $\beta_1 - \beta = 0$  in  $H_{p-1}(C_\bullet)$ . This shows that  $\bar{\beta}$  is indeed independent of the choice of  $\alpha'$ .

We now show that  $\text{Im}(H_p(\partial)) = \text{Ker}(H_{p-1}(\phi_\bullet))$ . Exactness at the other terms is easy to verify and is left as an exercise.

Let  $\beta \in C_{p-1}$  be a cycle such that  $\bar{\beta} \in H_{p-1}(C_\bullet)$  is in the image of  $H_p(\partial)$ . Let  $\alpha'' \in C''_p$  be such that  $H_p(\partial)(\bar{\alpha''}) = \bar{\beta}$  and let  $\alpha' \in C'_p, \beta' \in C'_{p-1}$  be as above. Then,  $\beta' = \phi_{p-1}(\beta) = \partial'_p(\alpha') \in B_{p-1}(C'_\bullet)$ . Descending to homology, this shows that  $H_{p-1}(\phi_\bullet)(\bar{\beta}) = 0$ , and  $\beta \in \text{Ker}(H_{p-1}(\phi_\bullet))$ .

Now, let  $\beta \in C_{p-1}$ , such that  $\bar{\beta} \in \text{Ker}(H_{p-1}(\phi_\bullet))$ . This implies that  $\phi_{p-1}(\beta) \in \text{Im}(\partial'_p)$ . Hence, there exists  $\alpha' \in C'_p$  such that  $\partial'_p(\alpha') = \phi_{p-1}(\beta)$ . Let  $\alpha'' = \psi_p(\alpha')$ . Since,  $\psi_{p-1}(\partial'_p(\alpha')) = \psi_{p-1}(\phi_{p-1}(\beta)) = 0$  by commutativity of the diagram, we have that  $\partial''_p(\alpha) = 0$ . Hence,  $\alpha$  is a cycle and it is easy to verify that  $H_p(\partial)(\bar{\alpha}) = \bar{\beta}$  and hence  $\bar{\beta} \in \text{Im}(H_p(\partial))$ .  $\square$

Another tool from homological algebra is the following Five Lemma.

**Lemma 6.11. [Five Lemma]** *Let*

$$\begin{array}{ccccccccc}
 C_1 & \xrightarrow{\phi_1} & C_2 & \xrightarrow{\phi_2} & C_3 & \xrightarrow{\phi_3} & C_4 & \xrightarrow{\phi_4} & C_5 \\
 \downarrow a & & \downarrow b & & \downarrow c & & \downarrow d & & \downarrow e \\
 D_1 & \xrightarrow{\psi_1} & D_2 & \xrightarrow{\psi_2} & D_3 & \xrightarrow{\psi_3} & D_4 & \xrightarrow{\psi_4} & D_5
 \end{array}$$

*be a commutative diagram such that each row is exact. Then if  $a, b, d, e$  are isomorphisms, so is  $c$ .*

**Proof:** We first show that  $c$  is injective. Let  $c(x_3) = 0$  for some  $x_3 \in C_3$ . Then  $d \circ \phi_3(x_3) = 0 \Rightarrow \phi_3(x_3) = 0$ , because  $d$  is an isomorphism. Hence,  $x_3 \in \text{ker}(\phi_3) = \text{Im}(\phi_2)$ . Let  $x_2 \in C_2$  be such that  $x_3 = \phi_2(x_2)$ . But then,  $\psi_2 \circ b(x_2) = 0 \Rightarrow b(x_2) \in \text{ker}(\psi_2) = \text{Im}(\psi_1)$ . Let  $y_1 \in D_1$  be such that  $\psi_1(y_1) = b(x_2)$ . Since  $a$  is an isomorphism there exists  $x_1 \in C_1$  such that  $y_1 = a(x_1)$  and  $\psi_1 \circ a(x_1) = b(x_2) = b \circ \phi_1(x_1)$ . Since  $b$  is an isomorphism this implies that  $x_2 = \phi_1(x_1)$ , and thus  $x_3 = \phi_2 \circ \phi_1(x_1) = 0$ .

Next we show that  $c$  is surjective. Let  $y_3 \in D_3$ . Since  $d$  is surjective there exists  $x_4 \in C_4$  such that  $\psi_3(y_3) = d(x_4)$ . Now,  $\psi_4 \circ \psi_3(y_3) = 0 = e \circ \phi_4(x_4)$ . Since  $e$  is injective this implies that  $x_4 \in \ker(\phi_4) = \text{Im}(\phi_3)$ . Let  $x_3 \in C_3$  be such that  $x_4 = \phi_3(x_3)$ . Then,  $d \circ \phi_3(x_3) = \psi_3 \circ c(x_3) = \psi_3(y_3)$ . Hence,  $c(x_3) - y_3 \in \ker(\psi_3) = \text{Im}(\psi_2)$ . Let  $y_2 \in D_2$  be such that  $\psi_2(y_2) = c(x_3) - y_3$ . There exists  $x_2 \in C_2$  such that  $\psi_2 \circ b(x_2) = c(x_3) - y_3 = c \circ \phi_2(x_2)$ . But then,  $c(x_3 - \phi_2(x_2)) = y_3$  showing that  $c$  is surjective.  $\square$

We use Lemma 6.10 to prove the existence of the so called Mayer-Vietoris sequence.

**Theorem 6.12. [Mayer-Vietoris]** *Let  $K$  be a simplicial complex and let  $K_1, K_2$  be sub-complexes of  $K$ . Then there is an exact sequence*

$$\dots \rightarrow H_p(K_1) \oplus H_p(K_2) \rightarrow H_p(K_1 \cup K_2) \rightarrow H_{p-1}(K_1 \cap K_2) \rightarrow \dots$$

**Proof:** We define homomorphisms  $\phi_\bullet, \psi_\bullet$  so that the sequence

$$0 \rightarrow C_\bullet(K_1 \cap K_2) \xrightarrow{\phi_\bullet} C_\bullet(K_1) \oplus C_\bullet(K_2) \xrightarrow{\psi_\bullet} C_\bullet(K_1 \cup K_2) \rightarrow 0$$

is exact.

There are natural inclusion homomorphisms  $i_1: C_\bullet(K_1 \cap K_2) \rightarrow C_\bullet(K_1)$  and  $i_2: C_\bullet(K_1 \cap K_2) \rightarrow C_\bullet(K_2)$ , as well as  $j_1: C_\bullet(K_1) \rightarrow C_\bullet(K_1 \cup K_2)$  and  $j_2: C_\bullet(K_2) \rightarrow C_\bullet(K_1 \cup K_2)$ .

For  $c \in C_\bullet(K_1 \cap K_2)$ , we define  $\phi(c) = (i_1(c), -i_2(c))$ .

For  $(d, e) \in C_\bullet(K_1) \oplus C_\bullet(K_2)$ , we define  $\psi(d, e) = j_1(d) + j_2(e)$ .

It is an exercise to check that, with these choices of  $\phi_\bullet$  and  $\psi_\bullet$ , the sequence

$$0 \rightarrow C_\bullet(K_1 \cap K_2) \xrightarrow{\phi_\bullet} C_\bullet(K_1) \oplus C_\bullet(K_2) \xrightarrow{\psi_\bullet} C_\bullet(K_1 \cup K_2) \rightarrow 0$$

is exact. Now, apply Lemma 6.10 to complete the proof.  $\square$

### 6.1.5 Chain Homotopy

We identify a property that guarantees that two chain homomorphisms induce identical homomorphisms in homology. The property is that they are chain homotopic.

Given two chain complexes,  $C_\bullet = (C_p, \partial_p)$  and  $C'_\bullet = (C'_p, \partial'_p)$ , two chain homomorphisms  $\phi_\bullet, \psi_\bullet: C_\bullet \rightarrow C'_\bullet$  are **chain homotopic** (denoted  $\phi_\bullet \sim \psi_\bullet$ ) if there exists a sequence of homomorphisms,  $\gamma_p: C_p \rightarrow C'_{p+1}$  such that

$$\partial'_{p+1} \circ \gamma_p + \gamma_{p-1} \circ \partial_p = \phi_p - \psi_p \tag{6.5}$$

for all  $p$ . The collection  $\gamma_\bullet$  of the homomorphisms  $\gamma_p$  is called a **chain homotopy** between  $C_\bullet$  and  $C'_\bullet$ .

$$\begin{array}{ccccccc}
\cdots & \xrightarrow{\partial_{p+2}} & C_{p+1} & \xrightarrow{\partial_{p+1}} & C_p & \xrightarrow{\partial_p} & \cdots \\
& & \searrow^{\gamma_{p+1}} & \downarrow^{\phi_{p+1}} & \swarrow^{\gamma_p} & \downarrow^{\phi_p} & \searrow^{\gamma_{p-1}} \\
& & & \psi_{p+1} & & \psi_p & \\
& & \swarrow_{\partial'_{p+2}} & \downarrow^{\psi_{p+1}} & \swarrow_{\partial'_{p+1}} & \downarrow^{\psi_p} & \swarrow_{\partial'_p} \\
\cdots & \xrightarrow{\partial'_{p+2}} & C'_{p+1} & \xrightarrow{\partial'_{p+1}} & C'_p & \xrightarrow{\partial'_p} & \cdots
\end{array}$$

**Lemma 6.13.** *Chain homotopy is an equivalence relation among chain homomorphisms from  $C_\bullet$  to  $C'_\bullet$ .*

**Proof:** Clearly every chain homomorphism  $\phi_\bullet: C_\bullet \rightarrow C'_\bullet$  is chain homotopic to itself (choose  $\gamma_p = 0$ ).

Also, if  $\gamma_p: C_p \rightarrow C'_{p+1}$  gives a chain homotopy between chain homomorphisms  $\phi_\bullet$  and  $\psi_\bullet$ , then  $-\gamma_\bullet$  gives a chain homotopy between  $\psi_\bullet$  and  $\phi_\bullet$ .

Finally, let  $\gamma_p: C_p \rightarrow C'_{p+1}$  be a chain homotopy between the chain homomorphisms  $\phi_\bullet$  and  $\psi_\bullet$  and let  $\lambda_p: C_p \rightarrow C'_{p+1}$  be a chain homotopy between the chain homomorphisms  $\psi_\bullet$  and  $\eta_\bullet$ .

Then, the homomorphisms  $\gamma_p + \lambda_p$  give a chain homotopy between  $\phi_\bullet$  and  $\eta_\bullet$ . This is because

$$\begin{aligned}
& \partial'_{p+1} \circ (\gamma_p + \lambda_p) + (\gamma_{p-1} + \lambda_{p-1}) \circ \partial_p \\
&= \partial'_{p+1} \circ \gamma_p + \gamma_{p-1} \circ \partial_p + \partial'_{p+1} \circ \lambda_p + \gamma_{p-1} \circ \lambda_p \\
&= \phi_p - \psi_p + \psi_p - \eta_p \\
&= \phi_p - \eta_p.
\end{aligned}$$

□

**Proposition 6.14.** *If  $\phi_\bullet \sim \psi_\bullet: C_\bullet \rightarrow C'_\bullet$ , then*

$$H_\star(\phi_\bullet) = H_\star(\psi_\bullet): H_\star(C_\bullet) \rightarrow H_\star(C'_\bullet).$$

**Proof:** Let  $c$  be a  $p$ -cycle in  $C_\bullet$ , that is  $c \in \text{Ker}(\partial_p)$ . Since  $\phi_\bullet$  and  $\psi_\bullet$  are chain homotopic, there exists a sequence of homomorphisms  $\gamma_p: C_p \rightarrow C'_{p+1}$  satisfying equation (6.1).

Thus,

$$(\partial'_{p+1} \circ \gamma_p + \gamma_{p-1} \circ \partial_p)(c) = (\phi_p - \psi_p)(c).$$

Now, since  $c$  is a cycle,  $\partial_p(c) = 0$ , and moreover,  $\partial'_{p+1}(\gamma_p(c))$  is a boundary. Thus,  $(\phi_p - \psi_p)(c) = 0$  in  $H_p(C'_\bullet)$ . □

*Example 6.15.* As an example of chain homotopy, consider the simplicial complex  $K$  whose simplices are all the faces of a single simplex  $[a_0, a_1, \dots, a_k] \subset \mathbb{R}^k$ . Consider the chain homomorphisms  $C_\bullet(\phi)$  and  $C_\bullet(\psi)$  induced by the simplicial maps  $\phi = \text{Id}_K$  and  $\psi$  such that  $\psi(a_i) = a_0, 1 \leq i \leq k$ .

Then,  $C_\bullet(\phi)$  and  $C_\bullet(\psi)$  are chain homotopic by the chain homotopy  $\gamma$  defined by  $\gamma_p([a_{i_0}, \dots, a_{i_p}]) = [a_0, a_{i_0}, \dots, a_{i_p}]$  for  $p \geq 0$  and  $\gamma_p = 0$  otherwise.

Clearly for  $p > 0$ ,

$$\begin{aligned} & (\partial_{p+1} \circ \gamma_p + \gamma_{p-1} \circ \partial_p)([a_{i_0}, \dots, a_{i_p}]) \\ &= [a_{i_0}, \dots, a_{i_p}] - [a_0, \partial_p([a_{i_0}, \dots, a_{i_p}])] + [a_0, \partial_p([a_{i_0}, \dots, a_{i_p}])] \\ &= [a_{i_0}, \dots, a_{i_p}] \\ &= (\phi_p - \psi_p)([a_{i_0}, \dots, a_{i_p}]). \end{aligned}$$

For  $p = 0$ ,

$$\begin{aligned} (\partial_1 \circ \gamma_0 + \gamma_{-1} \circ \partial_0)([a_{i_0}]) &= [a_{i_0}] - [a_0] \\ &= (\phi_0 - \psi_0)([a_{i_0}]). \end{aligned}$$

It is now easy to deduce that  $H_0(K) = \mathbb{Q}$  and  $H_i(K) = 0$  for all  $i > 0$ . □

A simplicial complex  $K$  is **acyclic** if  $H_0(K) = \mathbb{Q}$ ,  $H_i(K) = 0$  for all  $i > 0$ .

**Lemma 6.16.** *Let  $K$  be the simplicial complex whose simplices are all the faces of a simplex  $s = [a_0, a_1, \dots, a_k] \subset \mathbb{R}^k$ . Then,  $K$  is acyclic, and its barycentric subdivision  $\text{ba}(K)$  is acyclic.*

**Proof:** Recall from Section 5.6 that for every ascending sequence of simplices in  $K$ ,

$$s_0 \prec s_1 \prec \dots \prec s_j,$$

the simplex  $[\text{ba}(s_0), \dots, \text{ba}(s_j)]$  is included in  $\text{ba}(K)$ .

Consider the chain homomorphisms,  $\phi, \psi: C_\bullet(\text{ba}(K)) \rightarrow C_\bullet(\text{ba}(K))$ , induced by the simplicial maps  $\text{Id}$  and  $\psi$  defined by  $\psi(\text{ba}(s_i)) = \text{ba}(s)$ , for each  $s_i \in K$ .

Then,  $\phi$  and  $\psi$  are chain homotopic by the chain homotopy  $\gamma$  defined by  $\gamma_p([\text{ba}(s_0), \dots, \text{ba}(s_p)]) = [\text{ba}(s), \text{ba}(s_0), \dots, \text{ba}(s_p)]$  for  $p \geq 0$  and  $\gamma_p = 0$  otherwise.

Clearly for  $p > 0$ ,

$$\begin{aligned} & (\partial_{p+1} \circ \gamma_p + \gamma_{p-1} \circ \partial_p)([\text{ba}(s_0), \dots, \text{ba}(s_p)]) \\ &= [\text{ba}(s_0), \dots, \text{ba}(s_p)] - [\text{ba}(s), \partial_p([\text{ba}(s_0), \dots, \text{ba}(s_p)])] \\ & \quad + [\text{ba}(s), \partial_p([\text{ba}(s_0), \dots, \text{ba}(s_p)])] \\ &= [\text{ba}(s_0), \dots, \text{ba}(s_p)] \\ &= (\phi_p - \psi_p)([\text{ba}(s_0), \dots, \text{ba}(s_p)]). \end{aligned}$$

For  $p = 0$ ,

$$\begin{aligned} (\partial_1 \circ \gamma_0 + \gamma_{-1} \circ \partial_0)([\text{ba}(s_i)]) &= [\text{ba}(s_i)] - [\text{ba}(s)] \\ &= (\phi_0 - \psi_0)([\text{ba}(s_i)]). \end{aligned}$$



It is now easy to deduce that  $H_0(\text{ba}(K)) = \mathbb{Q}$  and  $H_i(\text{ba}(K)) = 0$  for all  $i > 0$ .  $\square$

We now identify a criterion that is sufficient to show that two homomorphisms are chain homotopic in the special case of chain complexes coming from simplicial complexes. The key notion is that of an **acyclic carrier function**.

Let  $K, L$  be two complexes. A function  $\xi$  which maps every simplex  $s \in K$  to a sub-complex  $\xi(s)$  of  $L$  is called a **carrier function** provided

$$s' \prec s \Rightarrow \xi(s') \subset \xi(s)$$

for all  $s, s' \in K$ . Moreover, if  $\xi(s)$  is acyclic for all  $s \in K$ ,  $\xi$  is called an **acyclic carrier function**. A chain homomorphism  $\phi_\bullet: C_\bullet(K) \rightarrow C_\bullet(L)$  is **carried by a carrier function**  $\xi$  if for all  $p$  and each  $s \in K_p$ ,  $\phi_p(s)$  is a chain in  $\xi(s)$ .

The most important property of a carrier function is the following.

**Lemma 6.17.** *If  $\phi_\bullet, \psi_\bullet: C_\bullet(K) \rightarrow C_\bullet(L)$  are chain homomorphisms carried by the same acyclic carrier  $\xi$ , then  $\phi_\bullet \sim \psi_\bullet$ .*

**Proof:** Let  $\partial$  (resp.  $\partial'$ ) be the boundary maps of  $C_\bullet(K)$  (resp.  $C_\bullet(L)$ ). We construct a chain homotopy  $\gamma$  dimension by dimension.

For  $s_0 \in C_0(K)$ ,  $\phi_0(s_0) - \psi_0(s_0)$  is a chain in  $\xi(s_0)$  which is acyclic. Since  $\xi(s_0)$  is acyclic,  $\phi_0(s_0) - \psi_0(s_0)$  must also be a boundary. Thus, there exists a chain  $t \in C_1(L)$  such that  $\partial'_1(t) = \phi_0(s_0) - \psi_0(s_0)$ , and we let  $\gamma_0(s_0) = t$ .

Now, assume that for all  $q < p$  we have constructed  $\gamma_q$  such that  $(\phi - \psi)_q = \partial'_{q+1} \circ \gamma_q + \gamma_{q-1} \circ \partial_q$  and  $\gamma_q(s) \subset \xi(s)$  for all  $q$ -simplices  $s$ .

We define  $\gamma_p(s)$  for  $p$ -simplices  $s$  and extend it linearly to  $p$ -chains. Notice first that  $(\phi - \psi)_p(s) \subset \xi(s)$  by hypothesis and that  $\gamma_{p-1} \circ \partial_p(s)$  is a chain in  $\xi(s)$  by the induction hypothesis. Hence  $((\phi - \psi)_p - \gamma_{p-1} \circ \partial_p)(s)$  is a chain in  $\xi(s)$  and let this chain be  $t$ . Then,

$$\begin{aligned} \partial'_p(t) &= \partial'_p((\phi - \psi)_p - \gamma_{p-1} \circ \partial_p)(s) \\ &= (\partial'_p \circ (\phi - \psi)_p - \partial'_p \circ \gamma_{p-1} \circ \partial_p)(s) \\ &= ((\phi - \psi)_{p-1} \circ \partial_p)(s) - (\partial'_p \circ \gamma_{p-1} \circ \partial_p)(s) \\ &= ((\phi - \psi)_{p-1} - \partial'_p \circ \gamma_{p-1})(\partial_p(s)) \\ &= \gamma_{p-2} \circ \partial_{p-1} \circ \partial_p(s) \\ &= 0 \end{aligned}$$

so that  $t$  is a cycle.

But, since  $t = ((\phi - \psi)_p - \gamma_{p-1} \circ \partial_p)(s)$  is a chain in  $\xi(s)$  and  $\xi(s)$  is acyclic,  $t$  must be a boundary as well. Thus, there is a chain,  $t'$ , such that  $t = \partial_{p+1}(t')$  and we define  $\gamma_p(s) = t'$ . It is straightforward to check that this satisfies all the conditions.  $\square$

Two simplicial maps  $\phi, \psi: K \rightarrow L$  are **contiguous** if  $\phi(s)$  and  $\psi(s)$  are faces of the same simplex in  $L$  for every simplex  $s \in K$ .

Two simplicial maps  $\phi, \psi: K \rightarrow L$  belong to the same contiguity class if there is a sequence of simplicial maps  $\phi_i, i = 0, \dots, n$ , such that  $\phi_0 = \phi, \phi_n = \psi$ , and  $\phi_i$  and  $\phi_{i+1}$  are contiguous for  $0 \leq i < n$ .

**Proposition 6.18.** *If the chain homomorphisms*

$$C_\bullet(\phi), C_\bullet(\psi): C_\bullet(K) \rightarrow C_\bullet(L)$$

*are induced by simplicial maps that belong to the same contiguity class, then  $H_\star(\phi) = H_\star(\psi)$ .*

**Proof:** We show that two contiguous simplicial maps induce chain homotopic chain homomorphisms, which will prove the proposition. In order to show this, we construct an acyclic carrier,  $\xi$ , for both  $\phi$  and  $\psi$ . For a simplex  $s \in K$ , let  $t$  be the smallest dimensional simplex of  $L$  such that  $\phi(s)$  and  $\psi(s)$  are both faces of  $t$  (in fact any such  $t$  will do). Let  $\xi(s)$  be the sub-complex of  $L$  consisting of all faces of  $t$ . Clearly,  $\xi$  is an acyclic carrier of both  $\phi$  and  $\psi$ , which implies that they are chain homotopic. □

### 6.1.6 The Simplicial Homology Groups Are Invariant Under Homeomorphism

We shall show that if  $K$  and  $L$  are two simplicial complexes that are homeomorphic then the homology groups of  $K$  and  $L$  are isomorphic.

#### 6.1.6.1 Homology and Barycentric Subdivision

The first step is to show that the homology groups of a simplicial complex  $K$  are isomorphic to those of its barycentric subdivision (see Definition page 182).

**Theorem 6.19.** *Let  $K$  be a simplicial complex and  $\text{ba}(K)$  its barycentric subdivision. Then,  $H_\star(K) \cong H_\star(\text{ba}(K))$ .*

As a consequence, we can iterate the operation of barycentric subdivision by setting  $K^{(1)} = \text{ba}(K)$  and, in general,  $K^{(n)} = \text{ba}(K^{(n-1)})$ , thus obtaining finer and finer subdivisions of the complex  $K$ .

**Corollary 6.20.** *Let  $K$  be a simplicial complex. Then,  $H_\star(K) \cong H_\star(K^{(n)})$  for all  $n > 0$ .*

In order to prove Theorem 6.19, we define some simplicial maps between the barycentric subdivision of a simplicial complex and the simplicial complex itself that will allow us to relate their homology groups.

Given a simplicial complex  $K$  and its barycentric subdivision  $\text{ba}(K)$ , a **Sperner map** is a map  $\omega: \text{ba}(K)_0 \rightarrow K_0$  such that  $\omega(\text{ba}(s))$  is one of the vertices of  $s$  for each simplex  $s \in K$ .

**Lemma 6.21.** *Any Sperner map can be linearly extended to a simplicial map.*

**Proof:** Let  $\omega: \text{ba}(K)_0 \rightarrow K_0$  be a Sperner map. Then an oriented simplex  $[\text{ba}(s_0), \dots, \text{ba}(s_i)]$  in  $\text{ba}(K)$  corresponds to  $s_0 < \dots < s_i$  in  $K$ , with  $\omega(\text{ba}(s_j)) \in s_j$ ,  $0 \leq j \leq i$ , and hence  $[\omega(\text{ba}(s_0)), \dots, \omega(\text{ba}(s_i))]$  is an oriented simplex in  $K$ .  $\square$

Given two simplicial complexes  $K, L$  and a simplicial map  $\phi: K \rightarrow L$ , there is a natural way to define a simplicial map  $\phi': \text{ba}(K) \rightarrow L'$  by setting  $\phi'(\text{ba}(s)) = \text{ba}(\phi(s))$  for every  $s \in K$  and extending it linearly to  $\text{ba}(K)$ . One can check that  $\phi'$  so defined is simplicial.

We define a new homomorphism

$$\alpha_\bullet: C_\bullet(K) \rightarrow C_\bullet(\text{ba}(K)),$$

which will play the role of an inverse to any Sperner map  $\omega$ , as follows:

It is defined on simplices recursively by,

$$\begin{aligned} \alpha_0(s) &= s, \\ \alpha_p(s) &= [\text{ba}(s), \alpha_{p-1}(\partial_p(s))], p > 0, \end{aligned}$$

(see Definition 6.1.1.1) and is then extended linearly to  $C_\bullet(K)$ . It is easy to verify that  $\alpha_\bullet$  is also a chain homomorphism.

**Lemma 6.22.** *Given a simplicial complex  $K$  and a Sperner map  $\omega$ ,  $C_\bullet(\omega) \circ \alpha_\bullet = \text{Id}_{C_\bullet(K)}$ .*

**Proof:** The proof is by induction on the dimension  $p$ . It is easily seen that the lemma holds when  $p=0$ . Consider a simplex  $s = [a_0, \dots, a_p]$ . Now,

$$\begin{aligned} (C_p(\omega) \circ \alpha_p)(s) &= C_p(\omega)([\text{ba}(s), \alpha_{p-1}(\partial_p(s))]) \\ &= [C_0(\omega)(\text{ba}(s)), C_{p-1}(\omega) \circ \alpha_{p-1}(\partial_p(s))]. \end{aligned}$$

By induction hypothesis,  $(C_{p-1}(\omega) \circ \alpha_{p-1})(\partial_p(s)) = \partial_p(s)$ . Since,  $C_0(\omega)(\text{ba}(s))$  is a vertex of  $s$  it follows that  $C_p(\omega) \circ \alpha_p(s) = s$ . This completes the induction.  $\square$

We now prove that  $\alpha_\bullet \circ C_\bullet(\omega) \sim \text{Id}_{C_\bullet(\text{ba}(K))}$  for a Sperner map  $\omega$ .

**Lemma 6.23.** *Let  $\omega: \text{ba}(K) \rightarrow K$  be a Sperner map. Then*

$$\alpha_\bullet \circ C_\bullet(\omega) \sim \text{Id}_{C_\bullet(\text{ba}(K))}.$$

**Proof:** We construct an acyclic carrier carrying  $\alpha_\bullet \circ C_\bullet(\omega)$  and  $\text{Id}_{C_\bullet(\text{ba}(K))}$ . Let  $\text{ba}(s)$  be a simplex of  $\text{ba}(K)$  and let  $b = \text{ba}(s)$  be the leading vertex of  $\text{ba}(s)$ , where  $s$  is a simplex in  $K$ . Let  $\xi(\text{ba}(s))$  be the sub-complex of  $\text{ba}(K)$  consisting of all simplices in the barycentric subdivision of  $s$ .

Clearly,  $\xi$  carries both  $\alpha_\bullet \circ C_\bullet(\omega)$  and  $\text{Id}_{C_\bullet(\text{ba}(K))}$  and is acyclic by Lemma 6.16, and hence satisfies the conditions for being an acyclic carrier.  $\square$

**Proof of Theorem 6.19:** Follows immediately from the preceding lemmas.  $\square$

### 6.1.6.2 Homeomorphisms Preserve Homology

Our goal in this paragraph is to show that homeomorphic polyhedra in real affine space have isomorphic homology groups.

**Theorem 6.24.** *If two simplicial complexes  $K \subset \mathbb{R}^k, L \subset \mathbb{R}^\ell$  are two simplicial complexes and  $f: |K| \rightarrow |L|$  is a homeomorphism, then there exists an isomorphism  $H_\star(f): H_\star(K) \rightarrow H_\star(L)$ .*

We will use the fact that our ground field is  $\mathbb{R}$  in two ways. In the next lemma, we use the fact that  $\mathbb{R}$  is (sequentially) compact in its metric topology in order to show the existence of a Lebesgue number for any finite open covering of a compact set in  $\mathbb{R}^k$ . Secondly, we will use the archimedean property of  $\mathbb{R}$ .

We first need a notation. For a vertex  $a$  of a simplicial complex  $K$ , its **star**  $\text{star}(a) \subset |K|$  is the union of the relative interiors of all simplices having  $a$  as a vertex, i.e.  $\text{star}(a) = \cup_{\{a\} \prec s} s^\circ$ . If the simplicial complexes  $K$  and  $L$  have the same polyhedron and if to every vertex  $a$  of  $K$  there is a vertex  $b$  of  $L$  such that  $\text{star}(a) \subset \text{star}(b)$ , then we write  $K < L$  and say  $K$  is **finer** than  $L$ .

It is clear that for any simplicial complex  $K$ ,  $K^{(n)} < K$ . Also, if  $K < L$  and  $L < M$  then  $K < M$ .

In the next lemma, we show that given a family of open sets whose union contains a compact subset  $S$  of  $\mathbb{R}^n$ , any "sufficiently small" subset of  $S$  is contained in a single set of the family.

We define the **diameter**  $\text{diam}(S)$  of a set  $S$  as the smallest number  $d$  such that  $S$  is contained in a ball of radius  $d/2$ .

**Lemma 6.25.** *Let  $\mathcal{A}$  be an open cover of a compact subset  $S$  of  $\mathbb{R}^n$ . Then, there exists  $\delta > 0$  (called the Lebesgue number of the cover) such that for any subset  $B$  of  $S$  with  $\text{diam}(B) < \delta$ , there exists an  $A \in \mathcal{A}$  such that  $B \subset A$ .*

**Proof:** Assume not. Then there exists a sequence of numbers  $\{\delta_n\}$  and sets  $S_n \subset S$  such that  $\delta_n \rightarrow 0$ ,  $\text{diam}(S_n) < \delta_n$ , and  $S_n \not\subset A$ , for all  $A \in \mathcal{A}$ .

Choose a point  $p_n$  in each  $S_n$ . Since  $S$  is compact, the sequence  $\{p_n\}$  has a convergent subsequence, and we pass to this subsequence and henceforth assume that the sequence  $\{p_n\}$  is convergent and its limit point is  $p$ .

Now  $p \in S$  since  $S$  is closed, and thus there exists a set  $A$  in the covering  $\mathcal{A}$  such that  $p \in A$ . Also, because  $A$  is open, there exists an  $\epsilon > 0$  such that the open ball  $B(p, \epsilon) \subset A$ .

Now choose  $n$  large enough so that  $\|p - p_n\| < \epsilon/2$  and  $\delta_n < \epsilon/2$ . We claim that  $S_n \subset A$ , which is a contradiction. To see this, observe that  $S_n$  contains a point  $p_n$  which is within  $\epsilon/2$  of  $p$ , but  $S_n$  also has diameter less than  $\epsilon/2$ . Hence it must be contained inside the ball  $B(p, \epsilon)$  and hence is contained in  $A$ .  $\square$

The **mesh**  $\text{mesh}(K)$  of a complex  $K$  is defined by

$$\text{mesh}(K) = \max \{ \text{diam}(s) \mid s \in K \}.$$

The following lemma bounds the mesh of the barycentric subdivision of a simplicial complex in terms of the mesh of the simplicial complex itself.

**Lemma 6.26.** *Let  $K$  be a simplicial complex of dimension  $k$ . Then,*

$$\text{mesh}(\text{ba}(K)) \leq \frac{k}{k+1} \text{mesh}(K).$$

**Proof:** First note that  $\text{mesh}(K)$  (resp.  $\text{mesh}(\text{ba}(K))$ ) equals the length of the longest edge in  $K$  (resp.  $\text{ba}(K)$ ). This follows from the fact that the diameter of a simplex equals the length of its longest edge.

Let  $(\text{ba}(s), \text{ba}(s'))$  be an edge in  $\text{ba}(K)$ , where  $s \prec s'$  are simplices in  $K$ . Also, without loss of generality, let  $s = [a_0, \dots, a_p]$  and  $s' = [a_0, \dots, a_q]$ .

Now,

$$\begin{aligned} \text{ba}(s) - \text{ba}(s') &= \frac{1}{p+1} \sum_{0 \leq i \leq p} a_i - \frac{1}{q+1} \sum_{0 \leq i \leq q} a_i \\ &= \left( \frac{1}{p+1} - \frac{1}{q+1} \right) \sum_{0 \leq i \leq p} a_i - \frac{1}{q+1} \sum_{p+1 \leq i \leq q} a_i \\ &= \frac{q-p}{q+1} \left( \frac{1}{p+1} \sum_{0 \leq i \leq p} a_i - \frac{1}{q-p} \sum_{p+1 \leq i \leq q} a_i \right). \end{aligned}$$

The points  $1/(p+1) \sum_{0 \leq i \leq p} a_i$  and  $1/(q-p) \sum_{p+1 \leq i \leq q} a_i$  are both in  $s'$ .

Hence, we have

$$\| \text{ba}(s) - \text{ba}(s') \| \leq \frac{q-p}{q+1} \text{mesh}(K) \leq \frac{q}{q+1} \text{mesh}(K) \leq \frac{k}{k+1} \text{mesh}(K). \quad \square$$

**Lemma 6.27.** *For any two simplicial complexes  $K$  and  $L$  such that  $|K| = |L| \subset \mathbb{R}^k$ , there exists  $n > 0$  such that  $K^{(n)} < L$ .*

**Proof:** The sets  $\text{star}(a)$  for each vertex  $a \in L_0$  give an open cover of the polyhedron  $|L|$ . Since the polyhedron is compact, for any such open cover there exists, by Lemma 6.25, a  $\delta > 0$  such that any subset of the polyhedron of diameter  $< \delta$  is contained in an element of the open cover, that is in  $\text{star}(a)$  for some  $a \in L_0$ . Using the fact that  $\text{mesh}(\text{ba}(K)) \leq k/(k + 1) \text{mesh}(K)$  and hence  $\text{mesh}(K^{(n)}) \leq (k/(k + 1))^n \text{mesh}(K)$ , we can choose  $n$  large enough so that for each  $b \in K_0^{(n)}$ , the set  $\text{star}(b)$  having diameter  $< 2 \text{mesh}(K^{(n)})$  is contained in  $\text{star}(a)$  for some  $a \in L_0$ . □

**Lemma 6.28.** *Let  $K, L$  be two simplicial complexes with  $K < L$ , and such that  $|K| = |L| \subset \mathbb{R}^k$ . Then, there exists a well-defined isomorphism*

$$i(K, L): H_*(K) \rightarrow H_*(L),$$

*which respects composition. In other words, given another simplicial complex  $M$  with  $|M| = |L|$  and  $L < M$ ,*

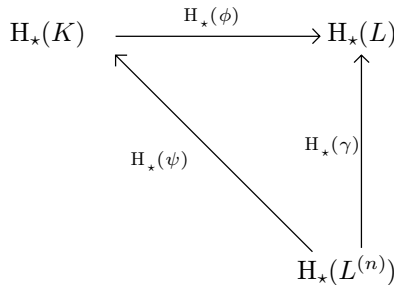
$$i(K, M) = i(L, M) \circ i(K, L).$$

**Proof:** Since  $K < L$ , for any vertex  $a \in K_0$ , there exists a vertex  $b \in L_0$  such that  $\text{star}(a) \subset \text{star}(b)$  since  $K < L$ . Consider a map  $\phi: K_0 \rightarrow L_0$  that sends each vertex  $a \in K_0$  to a vertex  $b \in L_0$  satisfying  $\text{star}(a) \subset \text{star}(b)$ . Notice that this agrees with the definition of a Sperner map in the case where  $K$  is a barycentric subdivision of  $L$ . Clearly, such a map is simplicial. Note that even though the simplicial map  $\phi$  is not uniquely defined, any other choice of the simplicial map satisfying the above condition is contiguous to  $\phi$  and thus induces the same homomorphism between  $H_*(K)$  and  $H_*(L)$ . Also, by Lemma 6.27, we can choose  $n$  such that  $L^{(n)} < K$  and a simplicial map  $\psi: L^{(n)} \rightarrow K$  that gives rise to a homomorphism

$$H_*(\psi): H_*(L^{(n)}) \rightarrow H_*(K).$$

In addition, using Theorem 6.19, we have an isomorphism

$$H_*(\gamma): H_*(L^{(n)}) \rightarrow H_*(L).$$



Again, note that the homomorphisms  $H_*(\psi), H_*(\gamma)$  are well-defined, even though the simplicial maps from which they are induced are not. Moreover,  $H_*(\gamma) = H_*(\phi) \circ H_*(\psi)$ . To see this let  $c \in L_0^{(n)}$ ,  $a = \psi(c) \in K_0$ , and  $b = \phi(a) \in L_0$ . Also, let  $b' = \gamma(c) \in L_0$ .

Then,  $\text{star}(c) \subset \text{star}(a)$  and  $\text{star}(a) \subset \text{star}(b)$ , so that  $\text{star}(c) \subset \text{star}(b)$ . Also,  $\text{star}(c) \subset \text{star}(b')$ .

Let  $s$  be the simplex in  $L$  such that  $b \in s$  and  $c \in s^\circ$ . Similarly, let  $t$  be the simplex in  $L$  such that  $b' \in t$  and  $c \in t^\circ$ . But, this implies that  $s^\circ \cap t^\circ \neq \emptyset$ , implying that  $s = t$ . This proves that the simplicial maps  $\phi \circ \psi$  and  $\gamma$  take a simplex  $s$  in  $L^{(n)}$  to faces of the simplex in  $L$  containing  $s$ , and hence,  $\phi \circ \psi$  and  $\gamma$  are contiguous, implying that  $H_*(\gamma) = H_*(\phi) \circ H_*(\psi)$ .

Now, since  $H_*(\gamma)$  is surjective, so is  $H_*(\phi)$ . The same reasoning for the pair  $L^{(n)} < K$  tells us that  $H_*(\psi)$  is surjective. Now, since  $H_*(\gamma)$  is injective and  $H_*(\psi)$  is surjective,  $H_*(\phi)$  is injective.

Define,  $i(K, L) = H_*(\phi)$ . Clearly,  $i(K, L)$  is independent of the particular simplicial map  $\phi$  chosen to define it. It also follows from the definition that the homomorphisms  $i$  respect composition.  $\square$

We next show that any continuous map between two polyhedrons can be suitably approximated by a simplicial map between some subdivisions of the two polyhedrons.

Given two simplicial complexes  $K, L$  and a continuous map  $f: |K| \rightarrow |L|$ , a simplicial map  $\phi: K \rightarrow L$  is a **simplicial approximation** to  $f$  if  $f(x) \in s^\circ$  implies  $\phi(x) \in s$ .

**Proposition 6.29.** *Given two simplicial complexes  $K, L$  and a continuous map  $f: |K| \rightarrow |L|$ , there exists an integer  $n > 0$  and a simplicial map  $\phi: K^{(n)} \rightarrow L$  that is a simplicial approximation to  $f$ .*

**Proof:** The family of open sets  $\{\text{star}(b) | b \in L_0\}$  is an open cover of  $L$ , and by continuity of  $f$  the family,  $\{f^{-1}(\text{star}(b)) | b \in L_0\}$  is an open cover of  $|K|$ . Let  $\delta$  be the Lebesgue number of this cover of  $|K|$  and choose  $n$  large enough so that  $\mu(K^{(n)}) < \delta/2$ . Thus, for every vertex  $a$  of  $K_0^{(n)}$ ,  $f(\text{star}(a)) \subset \text{star}(b)$  for some  $b \in L_0$ . It is easy to see that the map which sends  $a$  to such a  $b$  for every vertex  $a \in K_0^{(n)}$  induces a simplicial map  $\phi: K^{(n)} \rightarrow L$ . To see this, let  $s = [a_0, \dots, a_m]$  be a simplex in  $K^{(n)}$ . Then, by the definition of  $\phi$ ,  $\bigcap_{i=0}^m \text{star}(\phi(a_i)) \neq \emptyset$  since it contains  $f(s)$ . Hence,  $\{\phi(a_i) \mid 0 \leq i \leq m\}$  must span a simplex in  $L$ .

We now claim that  $\phi$  is a simplicial approximation to  $f$ . Let  $x \in |K|$  such that  $x \in s^\circ$  for a simplex  $s$  in  $K^{(n)}$ , and let  $f(x) \in t^\circ \subset |L|$ .

Let  $a \in K_0^{(n)}$  be a vertex of  $s$ , and let  $b = \phi(a)$ . From the definition of  $\phi$ , we have that  $f(\text{star}(a)) \subset \text{star}(b)$ , and since  $x \in \text{star}(a)$ ,  $f(x) \in \text{star}(b)$ . Thus,  $f(x) \in \bigcap_{a \in s} \text{star}(\phi(a))$ , and hence each  $\phi(a)$ ,  $a \in s$ , is a vertex of the simplex  $t$ . Moreover, since  $\phi(x)$  lies in the simplex spanned by  $\{\phi(a) \mid a \in s\}$ , it is clear that  $\phi(x) \in t$ .  $\square$

**Proposition 6.30.** *Given any two simplicial complexes  $K$  and  $L$ , as well as a continuous map  $f: |K| \rightarrow |L|$ , there exists a well-defined homomorphism  $H_*(f): H_*(K) \rightarrow H_*(L)$  such that if  $N$  is another simplicial complex and  $g: |L| \rightarrow |N|$  is a continuous map, then*

$$H_*(g \circ f) = H_*(g) \circ H_*(f) \text{ and } H_*(\text{Id}_{|K|}) = \text{Id}_{H_*(K)}.$$

**Proof:** Choose  $n_1$  large enough so that there is a simplicial approximation  $\phi: K^{(n_1)} \rightarrow L$  to  $f$ . Define  $H_*(f) = H_*(\phi) \circ i(K^{(n_1)}, K)^{-1}$ . It is easy using Lemma 6.28 to see that  $H_*(f)$  does not depend on the choice of  $n_1$ .

Now, suppose that we choose simplicial approximations  $\psi: L^{(n_2)} \rightarrow N$  of  $g$  and  $\phi: K^{(n_1)} \rightarrow L^{(n_2)}$  of  $f$ .

$$\begin{aligned} i(L^{(n_2)}, L) \circ H_*(\phi) \circ i(K^{(n_1)}, K)^{-1} &= H_*(f), \\ H_*(\psi) \circ i(L^{(n_2)}, L)^{-1} &= H_*(g), \\ H_*(g) \circ H_*(f) &= H_*(\psi) \circ H_*(\phi) \circ i(K^{(n_1)}, K)^{-1}. \end{aligned}$$

Note that  $\psi \circ \phi: K^{(n_1)} \rightarrow N$  is a simplicial approximation of  $g \circ f$ . To see this, observe that for  $x \in |K|$ ,  $f(x) \in s$  implies that  $\phi(x) \in s$ , where  $s$  is a simplex in  $L^{(n_2)}$ . Since,  $\psi$  is a simplicial map  $\psi(f(x)) \in t$  implies that  $\psi(\phi(x)) \in t$  for any simplex  $t$  in  $N$ . This proves that  $\psi \circ \phi$  is a simplicial approximation of  $g \circ f$  and hence

$$H_*(\psi) \circ H_*(\phi) \circ i(K^{(n_1)}, K)^{-1} = H_*(g \circ f).$$

The remaining property that  $H_*(\text{Id}_{|K|}) = \text{Id}_{H_*(K)}$  is now easy to check. □

Theorem 6.24 is now an immediate consequence of Proposition 6.30.

### 6.1.6.3 Semi-algebraic Homeomorphisms Preserve Homology

We next prove a result similar to Theorem 6.24 for semi-algebraic homeomorphisms between polyhedra defined over any real closed field.

Let  $K$  and  $L$  be two simplicial complexes contained in  $\mathbb{R}^k$  whose vertices have rational coordinates. Since  $K$  and  $L$  have vertices with rational coordinates, they can be described by linear inequalities with rational coefficients and hence they are semi-algebraic subsets of  $\mathbb{R}^k$ . We denote by  $\text{Ext}(|K|, \mathbb{R})$  and  $\text{Ext}(|L|, \mathbb{R})$  the polyhedron defined by the same inequalities over  $\mathbb{R}$ .

**Theorem 6.31.** *Let  $K$  and  $L$  be two simplicial complexes whose vertices have rational coordinates. If  $\text{Ext}(|K|, \mathbb{R})$  and  $\text{Ext}(|L|, \mathbb{R})$  are semi-algebraically homeomorphic for a real closed field  $\mathbb{R}$ , then  $H_*(K) \cong H_*(L)$ .*

The theorem will follow from the transfer property stated in the next lemma.



**Lemma 6.32.** *Let  $K$  and  $L$  be two simplicial complexes whose vertices have rational coordinates. The following are equivalent*

- *There exists a semi-algebraic homeomorphism from  $\text{Ext}(|K|, \mathbb{R}_{\text{alg}})$  to  $\text{Ext}(|L|, \mathbb{R}_{\text{alg}})$ .*
- *There exists a semi-algebraic homeomorphism from  $\text{Ext}(|K|, \mathbb{R})$  to  $\text{Ext}(|L|, \mathbb{R})$  for a real closed field  $\mathbb{R}$ .*

**Proof:** It is clear that if  $g: |K| = \text{Ext}(|K|, \mathbb{R}_{\text{alg}}) \rightarrow |L| = \text{Ext}(|L|, \mathbb{R}_{\text{alg}})$  is a semi-algebraic homeomorphism, then  $\text{Ext}(g, \mathbb{R}): \text{Ext}(|K|, \mathbb{R}) \rightarrow \text{Ext}(|L|, \mathbb{R})$  is a semi-algebraic homeomorphism, using the properties of the extension stated in Chapter 2 Exercise 2.16), since the property of a semi-algebraic function  $g$  of being a semi-algebraic homeomorphism can be described by a formula.

Conversely let  $\mathbb{R}$  be a real closed field, and let  $f: \text{Ext}(|K|, \mathbb{R}) \rightarrow \text{Ext}(|L|, \mathbb{R})$  be a semi-algebraic homeomorphism. Let  $A = (a_1, \dots, a_N) \in \mathbb{R}^N$  be the vector of all the constants appearing in the definition of the semi-algebraic maps  $f$ .

Let  $\Gamma_f \subset \mathbb{R}^{2k}$  denote the graph of the semi-algebraic map  $f$ , and let  $\phi_f(Z_1, \dots, Z_{2k})$  denote the formula defining  $\Gamma$ . For  $1 \leq i \leq N$ , replace every appearance of the constant  $a_i$  in  $\phi_f$  by a new variable  $Y_i$  to obtain a new formula  $\psi$  with  $N + 2k$  variables  $Y_1, \dots, Y_N, Z_1, \dots, Z_{2k}$ . All constants appearing in  $\psi$  are now rational numbers.

For  $b \in \mathbb{R}^N$ , let  $\Gamma_f(b) \subset \mathbb{R}^{2k}$  denote the set defined by  $\psi(b, Z_1, \dots, Z_{2k})$ .

We claim that we can write a formula  $\Phi_f(Y_1, \dots, Y_N)$  such that, for every  $b \in \mathbb{R}^N$  satisfying  $\Phi_f$ , the set  $\Gamma_f(b) \subset \mathbb{R}^{2k}$  is the graph of a semi-algebraic homeomorphism from  $\text{Ext}(|K|, \mathbb{R})$  to  $\text{Ext}(|L|, \mathbb{R})$  (with the domain and range corresponding to the first  $k$  and last  $k$  coordinates respectively).

A semi-algebraic homeomorphism is a continuous, 1-1, and onto map, with a continuous inverse. Hence, in order to write such a formula, we first write formulas guaranteeing continuity, injectivity, surjectivity, and continuity of the inverse separately and then take their conjunction.

Thinking of  $\bar{Y} = (Y_1, \dots, Y_N)$  as parameters, let  $\Phi_1(\bar{Y})$  be the first-order formula expressing that given  $\bar{Y}$ , for every open ball  $B \subset \mathbb{R}^k$ , the set in  $\mathbb{R}^k$  defined by

$$\{(Z_1, \dots, Z_k) \mid \exists((Z_{k+1}, \dots, Z_{2k}) \in B \wedge \psi(\bar{Y}, Z_1, \dots, Z_{2k}))\}$$

is open in  $\mathbb{R}^k$ . Since, we can clearly quantify over all open balls in  $\mathbb{R}^k$  (quantify over all centers and radii), we can thus express the property of being open by a first-order formula,  $\Phi_1(\bar{Y})$ .

Similarly, it is an easy exercise to translate the properties of a semi-algebraic map being injective, surjective and having a continuous inverse, into formulas  $\Phi_2(\bar{Y})$ ,  $\Phi_3(\bar{Y})$ ,  $\Phi_4(\bar{Y})$ , respectively. Finally, to ensure that  $\Gamma_f(b)$  is the graph of a map from  $\text{Ext}(|K|, \mathbb{R})$  to  $\text{Ext}(|L|, \mathbb{R})$ , we recall that  $\text{Ext}(|K|, \mathbb{R})$  is defined by inequalities with rational coefficients and we can clearly write a formula  $\Phi_5(\bar{Y})$  having the required property.

Now, take  $\Phi_f = \Phi_1 \wedge \Phi_2 \wedge \Phi_3 \wedge \Phi_4 \wedge \Phi_5$ . Since we know that  $(a_1, \dots, a_N) \in R^N$  satisfies  $\Phi_f$ , and thus  $\exists Y_1, \dots, Y_N \Phi_f(Y_1, \dots, Y_N)$  is true in  $R$  by the Tarski-Seidenberg principle (see Theorem 2.80), it is also true over  $\mathbb{R}_{\text{alg}}$ . Hence, there exists  $(b_1, \dots, b_N) \in \mathbb{R}_{\text{alg}}^N$  that satisfies  $\Phi$ . By substituting  $(b_1, \dots, b_N)$  for  $(a_1, \dots, a_N)$  in the description of  $f$ , we obtain a description of a semi-algebraic homeomorphism

$$g: |K| = \text{Ext}(|K|, \mathbb{R}_{\text{alg}}) \rightarrow |L| = \text{Ext}(|L|, \mathbb{R}_{\text{alg}}). \quad \square$$

**Proof of Theorem 6.31:** Let  $f: \text{Ext}(|K|, \mathbb{R}) \rightarrow \text{Ext}(|L|, \mathbb{R})$  be a semi-algebraic homeomorphism. Using Lemma 6.32, there exists a semi-algebraic homeomorphism  $g: |K| = \text{Ext}(|K|, \mathbb{R}) \rightarrow |L| = \text{Ext}(|L|, \mathbb{R})$ . Hence,  $H_*(K)$  and  $H_*(L)$  are isomorphic using Theorem 6.24.  $\square$

## 6.2 Simplicial Homology of Closed and Bounded Semi-algebraic Sets

### 6.2.1 Definitions and First Properties

We first define the simplicial homology groups of a closed and bounded semi-algebraic set  $S$ .

By Theorem 5.43, a closed, bounded semi-algebraic set  $S$  can be triangulated by a simplicial complex  $K$  with rational coordinates. Choose a semi-algebraic triangulation  $f: |K| \rightarrow S$ . **The homology groups  $H_p(S)$  are  $H_p(K)$ ,  $p \geq 0$ .** We denote by  $H_*(S)$  the chain complex  $(H_p(S), 0)$  and call it the **homology of  $S$** .

That the homology  $H_*(S)$  does not depend on a particular triangulation up to isomorphism follows from the results of Section 6.1. Given any two triangulations,  $f: |K| \rightarrow S, g: |L| \rightarrow S$ , there exists a semi-algebraic homeomorphism,  $\phi = g^{-1}f: |K| \rightarrow |L|$ , and hence, using Theorem 6.31,  $H_*(K)$  and  $H_*(L)$  are isomorphic.

Note that two semi-algebraically homeomorphic closed and bounded semi-algebraic sets have isomorphic homology groups. Note too that the homology groups of  $S$  and those of its extension to a bigger real closed field are also isomorphic.

The homology groups of  $S$  are all finite dimensional vector spaces over  $\mathbb{Q}$  (see Definition 6.2). The dimension of  $H_p(S)$  as a vector space over  $\mathbb{Q}$  is called the  **$p$ -th Betti number** of  $S$  and denoted  $b_p(S)$ .

$$b(S) = \sum_i b_i(S)$$

the sum of the Betti numbers of  $S$ . The **Euler-Poincaré characteristic** of  $S$  is

$$\chi(S) = \sum_i (-1)^i b_i(S).$$

Note that  $\chi(\emptyset) = 0$ .

Using Proposition 6.3 and Theorem 5.43, we have the following result.

**Proposition 6.33.** *Let  $S \subset \mathbb{R}^k$  be a closed and bounded semi-algebraic set,  $K$  be a simplicial complex in  $\mathbb{R}^k$  and  $h: |K| \rightarrow S$  be a semi-algebraic homeomorphism. Let  $n_i(K)$  be the number of simplexes of dimension  $i$  of  $K$ . Then*

$$\chi(S) = \sum_i (-1)^i n_i(K).$$

In particular the Euler-Poincaré characteristic of a finite set of points is the cardinality of this set.

**Proposition 6.34.** *The number of connected components of a non-empty, closed, and bounded semi-algebraic set  $S$  is  $b_0(S)$ .*

**Proof:** Let  $f: |K| \rightarrow S$  be a triangulation of  $S$ . Hence,

$$H_0(S) \cong H_0(K).$$

Now apply Proposition 6.5. □

We now use Theorem 6.12 to relate the homology groups of the union and intersection of two closed and bounded semi-algebraic sets.

**Theorem 6.35. [Semi-algebraic Mayer-Vietoris]** *Let  $S_1, S_2$  be two closed and bounded semi-algebraic sets. Then there is an exact sequence*

$$\cdots \rightarrow H_p(S_1 \cap S_2) \rightarrow H_p(S_1) \oplus H_p(S_2) \rightarrow H_p(S_1 \cup S_2) \rightarrow H_{p-1}(S_1 \cap S_2) \rightarrow \cdots$$

**Proof:** We first obtain a triangulation of  $S_1 \cup S_2$  that is simultaneously a triangulation of  $S_1$ ,  $S_2$ , and  $S_1 \cap S_2$  using Theorem 5.43. We then apply Theorem 6.12. □

From the exactness of the Mayer-Vietoris sequence, we have the following corollary.

**Corollary 6.36.** *Let  $S_1, S_2$  be two closed and bounded semi-algebraic sets. Then,*

$$\begin{aligned} b_i(S_1) + b_i(S_2) &\leq b_i(S_1 \cup S_2) + b_i(S_1 \cap S_2), \\ b_i(S_1 \cap S_2) &\leq b_i(S_1) + b_i(S_2) + b_{i+1}(S_1 \cup S_2), \\ b_i(S_1 \cup S_2) &\leq b_i(S_1) + b_i(S_2) + b_{i-1}(S_1 \cap S_2), \\ \chi(S_1 \cup S_2) &= \chi(S_1) + \chi(S_2) - \chi(S_1 \cap S_2). \end{aligned}$$

**Proof:** Follows directly from Theorem 6.35. □

The Mayer-Vietoris sequence provides an easy way to compute the homology groups of some simple sets.

**Proposition 6.37.** *Consider the  $(k-1)$ -dimensional unit sphere  $S^{k-1} \subset \mathbb{R}^k$  for  $k > 1$ . If  $k \geq 0$ ,*

$$\begin{aligned} H_0(B_k) &\cong \mathbb{Q}, \\ H_i(B_k) &\cong 0, \quad i > 0. \end{aligned}$$

If  $k > 1$ ,

$$\begin{aligned} H_0(S^{k-1}) &\cong \mathbb{Q}, \\ H_i(S^{k-1}) &\cong 0, \quad 0 < i < k-1, \\ H_{k-1}(S^{k-1}) &\cong \mathbb{Q}, \\ H_i(S^{k-1}) &\cong 0, \quad i > k-1. \end{aligned}$$

**Proof:** We can decompose the unit sphere into two closed hemispheres,  $A, B$ , intersecting at the equator.

Each of the sets  $A, B$  is homeomorphic to the standard  $(k-1)$ -dimensional simplex, and  $A \cap B$  is homeomorphic to the  $(k-2)$ -dimensional sphere  $S^{k-2}$ .

If  $k = 1$ , it is clear that  $H_0(S^1) \cong H_1(S^1) \cong \mathbb{Q}$ .

For  $k \geq 2$ , the statement is proved by induction on  $k$ .

Assume that the result holds for spheres of dimensions less than  $k-1 \geq 2$ . The Mayer-Vietoris sequence for homology (Theorem 6.35) gives the exact sequence

$$\cdots \rightarrow H_p(A \cap B) \rightarrow H_p(A) \oplus H_p(B) \rightarrow H_p(A \cup B) \rightarrow H_{p-1}(A \cap B) \rightarrow \cdots$$

Here, the homology groups of  $A$  and  $B$  are isomorphic to those of a  $(k-1)$ -dimensional closed ball, and thus  $H_0(A) \cong H_0(B) \cong \mathbb{Q}$  and  $H_p(A) \cong H_p(B) \cong 0$  for all  $p > 0$ . Moreover, the homology groups of  $A \cap B$  are isomorphic to those of a  $(k-2)$ -dimensional sphere, and thus  $H_0(A \cap B) \cong H_{k-2}(A \cap B) \cong \mathbb{Q}$  and  $H_p(A \cap B) \cong 0$ , for  $p \neq k-2, p \neq 0$ . It now follows from the exactness of the above sequence that  $H_0(A \cup B) \cong H_{k-1}(A \cup B) \cong \mathbb{Q}$ , and  $H_p(A \cup B) \cong 0$ , for  $p \neq k-1, p \neq 0$ .

To see this, observe that the exactness of

$$H_{k-1}(A) \oplus H_{k-1}(B) \rightarrow H_{k-1}(A \cup B) \rightarrow H_{k-2}(A \cap B) \rightarrow H_{k-2}(A) \oplus H_{k-2}(B)$$

is equivalent to the following sequence being exact:

$$0 \rightarrow H_{k-1}(A \cup B) \rightarrow \mathbb{Q} \rightarrow 0,$$

and this implies that the homomorphism  $H_{k-1}(A \cup B) \rightarrow \mathbb{Q}$  is an isomorphism.  $\square$

## 6.2.2 Homotopy

Let  $X, Y$  be two topological spaces. Two continuous functions  $f, g: X \rightarrow Y$  are **homotopic** if there is a continuous function  $F: X \times [0, 1] \rightarrow Y$  such that  $F(x, 0) = f(x)$  and  $F(x, 1) = g(x)$  for all  $x \in X$ . Clearly, homotopy is an equivalence relation among continuous maps from  $X$  to  $Y$ . It is denoted by  $f \sim g$ .

The sets  $X, Y$  are **homotopy equivalent** if there exist continuous functions  $f: X \rightarrow Y$ ,  $g: Y \rightarrow X$  such that  $g \circ f \sim \text{Id}_X$ ,  $f \circ g \sim \text{Id}_Y$ . If two sets are homotopy equivalent, we also say that they have the same homotopy type.

Let  $X$  be a topological space and  $Y$  a closed subset of  $X$ . A **deformation retraction** from  $X$  to  $Y$  is a continuous function  $h: X \times [0, 1] \rightarrow X$  such that  $h(-, 0) = \text{Id}_X$  and such that  $h(-, 1)$  has its values in  $Y$  and such that for every  $t \in [0, 1]$  and every  $x$  in  $Y$ ,  $h(x, t) = x$ . If there is a deformation retraction from  $X$  to  $Y$ , then  $X$  and  $Y$  are clearly homotopy equivalent.

**Theorem 6.38.** *Let  $K, L$  be simplicial complexes over  $\mathbb{R}$  and  $f, g$  continuous homotopic maps from  $|K|$  to  $|L|$ . Then*

$$H_*(f) = H_*(g): H_*(K) \rightarrow H_*(L).$$

The proposition follows directly from the following two lemmas.

**Lemma 6.39.** *Let  $K, L$  be simplicial complexes over  $\mathbb{R}$  and  $\delta$  the Lebesgue number of the cover  $\{\text{star}(b) | b \in L_0\}$ . Let  $f, g: |K| \rightarrow |L|$  be two continuous maps. If  $\sup_{x \in |K|} |f(x) - g(x)| < \delta/3$ , then  $f$  and  $g$  have a common simplicial approximation.*

**Proof:** For  $b \in L_0$  let  $B_b = \{x \in |L| \mid \text{dist}(x, |L| - \text{star}(b)) > \delta/3\}$ . We first claim that  $b \in B_b$  and hence, the family of sets  $\{B_b | b \in L_0\}$  is an open covering of  $L$ . Consider the set  $|L| \cap B(b, 2\delta/3)$ . If  $|L| \cap B(b, 2\delta/3) \subset \text{star}(b)$ , then clearly,  $b \in B_b$ . Otherwise, since  $\text{diam}(|L| \cap B(b, 2\delta/3)) < \delta$ , there must exist a  $b' \in L_0$  such that  $|L| \cap B(b, 2\delta/3) \subset \text{star}(b')$ . But, then  $b \in \text{star}(b')$  implying that  $b = b'$ , which is a contradiction.

Let  $\epsilon$  be the Lebesgue number of the open cover of  $|K|$  given by  $\{f^{-1}(B_b) | b \in L_0\}$ . Then, there is an integer  $n$  such that  $\mu(K^{(n)}) < \epsilon/2$ . To every vertex  $a \in K^{(n)}$ , there is a vertex  $b \in L_0$  such that  $\text{star}(a) \subset f^{-1}(B_b)$ , and this induces a simplicial map,  $\phi: K^{(n)} \rightarrow L$ , sending  $a$  to  $b$ . We now claim that  $\phi$  is a simplicial approximation to both  $f$  and  $g$ .

Let  $x \in |K|$  such that  $x \in s^\circ$  for a simplex  $s$  in  $K^{(n)}$ , and let  $f(x) \in t_1^\circ$  and  $g(x) \in t_2^\circ$  for simplices  $t_1, t_2 \in L$ . Let  $a \in K_0^{(n)}$  be a vertex of  $s$ , and let  $b = \phi(a)$ . From the definition of  $\phi$ , we have that  $f(\text{star}(a)) \subset B_b \subset \text{star}(b)$ , and since  $x \in \text{star}(a)$ ,  $f(x) \in \text{star}(b)$ . Moreover, since  $|f(x) - g(x)| < \delta/3$  for all  $x \in |K|$ ,  $\text{dist}(f(\text{star}(a)), g(\text{star}(a))) < \delta/3$ , and hence  $g(\text{star}(a)) \subset \text{star}(b)$  and  $g(x) \in \text{star}(b)$ .

Thus,  $f(x) \in \cap_{a \in s} \text{star}(\phi(a))$ , and hence each  $\phi(a)$ ,  $a \in s$  is a vertex of the simplex  $t_1$ . Moreover, since  $\phi(x)$  lies in the simplex spanned by  $\{\phi(a) | a \in s\}$ , it is clear that  $\phi(x) \in t_1$ . Similarly,  $\phi(x) \in t_2$ , and hence  $\phi$  is simultaneously a simplicial approximation to both  $f$  and  $g$ .  $\square$

**Lemma 6.40.** *Let  $K, L$  be simplicial complexes over  $\mathbb{R}$  and suppose that  $f, g$  are homotopic maps from  $|K|$  to  $|L|$ . Then, there is an integer  $n$  and simplicial maps  $\phi, \psi: K^{(n)} \rightarrow L$  that are in the same contiguity class and such that  $\phi$  (resp.  $\psi$ ) is a simplicial approximation of  $f$  (resp.  $g$ ).*

**Proof:** Since  $f \sim g$ , there is a continuous map  $F: |K| \times [0, 1] \rightarrow |L|$  such that  $F(x, 0) = f(x)$  and  $F(x, 1) = g(x)$ . To a Lebesgue number  $\delta$  of the cover  $\text{star}(b)$ ,  $b \in L_0$ , there exists a number  $\epsilon$  such that  $|t - t'| < \epsilon$  implies  $\sup |F(x, t) - F(x, t')| < \delta/3$ . This follows from the uniform continuity of  $F$  since  $K \times [0, 1]$  is compact.

We now choose a sequence  $t_0 = 0 < t_1 < t_2 < \dots < t_n = 1$  such that  $|t_{i+1} - t_i| < \epsilon$  and let  $F(x, t_i) = f_i(x)$ . By the previous lemma,  $f_i$  and  $f_{i+1}$  have a common simplicial approximation  $\psi_i: K^{(n_i)} \rightarrow L$ . Let  $n = \max_i n_i$ , and let  $\phi_i: K^{(n)} \rightarrow L$  be the simplicial map induced by  $\psi_i$ . For each  $i$ ,  $0 \leq i < n$ ,  $\phi_i$  and  $\phi_{i+1}$  are contiguous and are simplicial approximations of  $f_i$  and  $f_{i+1}$  respectively. Moreover,  $\phi_0$  is a simplicial approximation of  $f$  and  $\phi_n$  a simplicial approximation of  $g$ . Hence, they are in the same contiguity class.  $\square$

We will now transfer the previous results to semi-algebraic sets and maps over a general real closed field  $\mathbb{R}$ . The method of transferring the results parallels those used at the end of Section 6.1.

Let  $X, Y$  be two closed and bounded semi-algebraic sets. Two semi-algebraic continuous functions  $f, g: X \rightarrow Y$  are **semi-algebraically homotopic**,  $f \sim_{sa} g$ , if there is a continuous semi-algebraic function  $F: X \times [0, 1] \rightarrow Y$  such that  $F(x, 0) = f(x)$  and  $F(x, 1) = g(x)$  for all  $x \in X$ . Clearly, semi-algebraic homotopy is an equivalence relation among semi-algebraic continuous maps from  $X$  to  $Y$ .

The sets  $X, Y$  are **semi-algebraically homotopy equivalent** if there exist semi-algebraic continuous functions  $f: X \rightarrow Y$ ,  $g: Y \rightarrow X$  such that  $g \circ f \sim_{sa} \text{Id}_X$ ,  $f \circ g \sim_{sa} \text{Id}_Y$ .

Let  $X$  be a closed and bounded semi-algebraic set and  $Y$  a closed semi-algebraic subset of  $X$ . A **semi-algebraic deformation retraction** from  $X$  to  $Y$  is a continuous semi algebraic function  $h: X \times [0, 1] \rightarrow X$  such that  $h(-, 0) = \text{Id}_X$  and such that  $h(-, 1)$  has its values in  $Y$  and such that for every  $t \in [0, 1]$  and every  $x$  in  $Y$ ,  $h(x, t) = x$ . If there is a semi-algebraic deformation retraction from  $X$  to  $Y$ , then  $X$  and  $Y$  are clearly semi-algebraically homotopy equivalent.

Using the transfer principle and the same technique used in the proof of Theorem 6.31, it is possible to prove,

**Proposition 6.41.** *Let  $K, L$  be simplicial complexes with rational vertices, and let  $f \sim_{sa} g$  be semi-algebraic continuous semi-algebraically homotopic maps from  $\text{Ext}(|K|, \mathbb{R})$  to  $\text{Ext}(|L|, \mathbb{R})$ . Then*

$$H_*(f) = H_*(g): H_*(K) \rightarrow H_*(L).$$

Finally, the following proposition holds in any real closed field.

**Theorem 6.42.** *Let  $\mathbb{R}$  be a real closed field. Let  $X, Y$  be two closed, bounded semi-algebraic sets of  $\mathbb{R}^k$  that are semi-algebraically homotopy equivalent. Then,  $H_*(X) \cong H_*(Y)$ .*

**Proof:** We first choose triangulations. Let  $\phi: |K| \rightarrow X$  and  $\psi: |L| \rightarrow Y$  be semi-algebraic triangulations of  $X$  and  $Y$ , respectively. Moreover, since  $X$  and  $Y$  are semi-algebraically homotopy equivalent, there exist semi-algebraic continuous functions  $f: X \rightarrow Y$ ,  $g: Y \rightarrow X$  such that  $g \circ f \sim_{sa} \text{Id}_X$ ,  $f \circ g \sim_{sa} \text{Id}_Y$ .

Then,  $f_1 = \psi^{-1} \circ f \circ \phi: |K| \rightarrow |L|$  and  $g_1 = \phi^{-1} \circ g \circ \psi: |L| \rightarrow |K|$  give a semi-algebraic homotopy equivalence between  $|K|$  and  $|L|$ . These are defined over  $\mathbb{R}$ . However, using the same method as in the proof of Lemma 6.32, we can show that in this case there exists  $f'_1: |K| \rightarrow |L|$  and  $g'_1: |L| \rightarrow |K|$  defined over  $\mathbb{R}$  giving a homotopy equivalence between  $|K|$  and  $|L|$ .

Now applying Proposition 6.41 and Proposition 6.30 we get

$$H_*(f'_1 \circ g'_1) = H_*(f'_1) \circ H_*(g'_1) = H_*(\text{Id}_K) = \text{Id}_{H_*(K)}$$

$$H_*(g'_1 \circ f'_1) = H_*(g'_1) \circ H_*(f'_1) = H_*(\text{Id}_L) = \text{Id}_{H_*(L)}.$$

This proves that  $H_*(X) = H_*(K) \cong H_*(L) = H_*(Y)$ .  $\square$

## 6.3 Homology of Certain Locally Closed Semi-Algebraic Sets

In Section 6.2 we have defined homology groups of closed and bounded semi-algebraic sets. Now, we consider more general semi-algebraic sets - namely, certain locally closed semi-algebraic sets. We first define homology groups for closed semi-algebraic sets, as well as for semi-algebraic sets which are realizations of sign conditions. These homology groups are homotopy invariant, but do not satisfy an additivity property useful in certain applications. In order to have a homology theory with the additivity property, we introduce the Borel-Moore homology groups and prove their basic properties.

### 6.3.1 Homology of Closed Semi-algebraic Sets and of Sign Conditions

We now define the homology groups for closed (but not necessarily bounded) semi-algebraic sets and for semi-algebraic sets defined by a single sign condition.

Let  $S \subset \mathbb{R}^k$  be any closed semi-algebraic set. By Proposition 5.49 (conic structure at infinity), there exists  $r \in \mathbb{R}$ ,  $r > 0$ , such that, for every  $r' \geq r$ , there exists a semi-algebraic deformation retraction from  $S$  to  $S_{r'} = S \cap B_k(0, r')$ , and there exists a semi-algebraic deformation retraction from  $S_r$  to  $S_{r'}$ . Thus the sets  $S_r$  and  $S_{r'}$  are homotopy equivalent. So, by Theorem 6.42,  $H(S_r) = H(S_{r'})$ .

**Notation 6.43. [Homology]** We define  $H_*(S) = H_*(S_r)$ .  $\square$

We have the following useful result.

**Proposition 6.44.** *Let  $S_1, S_2$  be two closed semi-algebraic sets. Then,*

$$\begin{aligned} b_i(S_1) + b_i(S_2) &\leq b_i(S_1 \cup S_2) + b_i(S_1 \cap S_2), \\ b_i(S_1 \cap S_2) &\leq b_i(S_1) + b_i(S_2) + b_{i+1}(S_1 \cup S_2), \\ b_i(S_1 \cup S_2) &\leq b_i(S_1) + b_i(S_2) + b_{i-1}(S_1 \cap S_2). \end{aligned}$$

**Proof:** Follows directly from Corollary 6.36 and the definition of the homology groups of a closed semi-algebraic set.  $\square$

We also define homology groups for semi-algebraic sets defined by a single sign condition.

Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{R}[X_1, \dots, X_k]$  be a set of  $s$  polynomials, and let  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$  be a realizable sign condition on  $\mathcal{P}$ . Without loss of generality, suppose

$$\begin{aligned} \sigma(P_i) &= 0 \quad \text{if } i = 1, \dots, j, \\ \sigma(P_i) &= 1 \quad \text{if } i = j + 1, \dots, \ell, \\ \sigma(P_i) &= -1 \quad \text{if } i = \ell + 1, \dots, s. \end{aligned}$$

We denote by  $\text{Reali}(\sigma) \subset \mathbb{R}^k$  the realization of  $\sigma$ . Let  $\delta > 0$  be a variable.

Consider the field  $\mathbb{R}\langle\delta\rangle$  of algebraic Puiseux series in  $\delta$ , in which  $\delta$  is an infinitesimal. Let  $\overline{\text{Reali}}(\sigma) \subset \mathbb{R}\langle\delta\rangle^k$  be defined by

$$\begin{aligned} \sum_{1 \leq i \leq k} X_i^2 \leq 1/\delta \wedge P_1 = \dots = P_j = 0 \\ \wedge P_{j+1} \geq \delta \wedge \dots \wedge P_\ell \geq \delta \wedge P_{\ell+1} \leq -\delta \wedge \dots \wedge P_s \leq -\delta. \end{aligned}$$

**Proposition 6.45.** *The set  $\overline{\text{Reali}}(\sigma)$  is a semi-algebraic deformation retract of the extension of  $\text{Reali}(\sigma)$  to  $\mathbb{R}\langle\delta\rangle$ .*

**Proof:** Consider the continuous semi-algebraic function  $f$  defined by

$$f(x) = \inf \left( 1, \frac{1}{X_1^2 + \dots + X_k^2}, \inf_{j+1 \leq i \leq s} (|P_i(x)|) \right)$$

and note that

$$\overline{\text{Reali}}(\sigma) = \{x \in \text{Ext}(\text{Reali}(\sigma), \mathbb{R}\langle\delta\rangle) \mid f(x) \geq \delta\}.$$

By Theorem 5.46 (Hardt's triviality), there exists  $t_0 \in \mathbb{R}$  such that

$$\begin{aligned} \{x \in \text{Reali}(\sigma) \mid t_0 \geq f(x) > 0\} \\ \text{(resp. } \{x \in \text{Ext}(\text{Reali}(\sigma), \mathbb{R}\langle\delta\rangle) \mid t_0 \geq f(x) \geq \delta\}) \end{aligned}$$

is homeomorphic to

$$\begin{aligned} \{x \in \text{Reali}(\sigma) \mid f(x) = t_0\} \times (0, t_0] \\ \text{(resp. } \{x \in \text{Ext}(\text{Reali}(\sigma), \mathbb{R}\langle\delta\rangle) \mid f(x) = t_0\} \times [\delta, t_0]). \end{aligned}$$



Moreover, the corresponding homeomorphisms  $\phi$  and  $\psi$  can be chosen such that  $\phi|_{\{x \in \text{Reali}(\sigma) \mid f(x) = t_0\}}$  and  $\psi|_{\{x \in \text{Ext}(\text{Reali}(\sigma, \mathbb{R}(\delta)) \mid f(x) = t_0\}}$  are identities.  $\square$

**Notation 6.46. [Homology of a sign condition]** We define

$$H_\star(\text{Reali}(\sigma)) = H_\star(\overline{\text{Reali}}(\sigma)). \quad \square$$

**Proposition 6.47.** *Suppose that  $\text{Reali}(\sigma)$  and  $\text{Reali}(\tau)$  are semi-algebraically homotopy equivalent, then*

$$H_\star(\text{Reali}(\sigma)) = H_\star(\text{Reali}(\tau)).$$

**Proof:** By Proposition 6.45,  $\overline{\text{Reali}}(\sigma)$  and  $\overline{\text{Reali}}(\tau)$  are homotopy equivalent. Now apply Theorem 6.42.  $\square$

**Exercise 6.1.** Consider the unit disk minus a point which is the set  $D$  defined by

$$X^2 + Y^2 - 1 < 0 \wedge X^2 + Y^2 > 0.$$

Prove that

$$\begin{aligned} H_0(D) &= \mathbb{Q}, \\ H_1(D) &= \mathbb{Q}, \\ H_2(D) &= 0. \end{aligned}$$

*Remark 6.48.* The homology groups we just defined agree with the singular homology groups [150] in the case when  $\mathbb{R} = \mathbb{R}$ : it is a consequence of Proposition 6.45 and the fact that the singular homology groups are homotopy invariants [150].  $\square$

### 6.3.2 Homology of a Pair

We now define the simplicial homology groups of pairs of closed and bounded semi-algebraic sets.

Let  $K$  be a simplicial complex and  $A$  a sub-complex of  $K$ . Then, there is a natural inclusion homomorphism,  $i: C_p(A) \rightarrow C_p(K)$ , between the corresponding chain groups. Defining the group  $C_p(K, A) = C_p(K)/i(C_p(A))$ , it is easy to see that the boundary maps  $\partial_p: C_p(K) \rightarrow C_{p-1}(K)$  descend to maps  $\partial_p: C_p(K, A) \rightarrow C_{p-1}(K, A)$ , so that we have a short exact sequence of complexes,

$$0 \rightarrow C_\bullet(A) \rightarrow C_\bullet(K) \rightarrow C_\bullet(K, A) \rightarrow 0.$$

Given a pair  $(K, A)$ , where  $A$  is a sub-complex of  $K$ , the group

$$H_p(K, A) = H_p(C_\bullet(K, A))$$

is the  **$p$ -th simplicial homology group** of the pair  $(K, A)$ .

It is clear from the definition that  $H_p(K, A)$  is a finite dimensional  $\mathbb{Q}$ -vector space. The dimension of  $H_p(K, A)$  as a  $\mathbb{Q}$ -vector space is called the  **$p$ -th Betti number of the pair**  $(K, A)$  and denoted  $b_p(K, A)$ . The **Euler-Poincaré characteristic** of the pair  $(K, A)$  is

$$\chi(K, A) = \sum_i (-1)^i b_i(K, A).$$

We now define the simplicial homology groups of a pair of closed and bounded semi-algebraic sets  $T \subset S \subset \mathbb{R}^k$ . By Theorem 5.43, such a pair of closed, bounded semi-algebraic sets can be triangulated using a pair of simplicial complexes  $(K, A)$  with rational coordinates, where  $A$  is a sub-complex of  $K$ . The  **$p$ -th simplicial homology group** of the pair  $(S, T)$ ,  $H_p(S, T)$ , is  $H_p(K, A)$ . The dimension of  $H_p(S, T)$  as a  $\mathbb{Q}$ -vector space is called the  **$p$ -th Betti number of the pair**  $(S, T)$  and denoted  $b_p(S, T)$ . The **Euler-Poincaré characteristic** of the pair  $(S, T)$  is

$$\chi(S, T) = \sum_i (-1)^i b_i(S, T).$$

**Exercise 6.2.** Consider the pair  $(S, T)$  where  $S$  is the closed unit disk defined by  $X^2 + Y^2 - 1 \leq 0$  and  $T$  is the union of the origin and the circle of radius one defined by  $X^2 + Y^2 - 1 = 0 \vee X^2 + Y^2 = 0$ .

Prove that

$$\begin{aligned} H_0(S, T) &= \mathbb{Q}, \\ H_1(S, T) &= \mathbb{Q}, \\ H_2(S, T) &= \mathbb{Q}. \end{aligned}$$

**Proposition 6.49.** Let  $T \subset S \subset \mathbb{R}^k$  be a pair of closed and bounded semi-algebraic set,  $(K, A)$  be a pair of simplicial complexes in  $\mathbb{R}^k$ , with  $A$  a sub-complex of  $K$  and let  $h: |K| \rightarrow S$  be a semi-algebraic homeomorphism such that the image of  $|K|$  is  $T$ . Then

$$\begin{aligned} \chi(S, T) &= \chi(K, A) \\ &= \chi(K) - \chi(A) \\ &= \chi(S) - \chi(T). \end{aligned}$$

**Proof:** From the short exact sequence of chain complexes,

$$0 \rightarrow C_\bullet(A) \rightarrow C_\bullet(K) \rightarrow C_\bullet(K, A) \rightarrow 0,$$

applying Lemma 6.10, we obtain the following long exact sequence of homology groups:

$$\cdots \rightarrow H_p(A) \rightarrow H_p(K) \rightarrow H_p(K, A) \rightarrow H_{p-1}(A) \rightarrow H_{p-1}(K) \rightarrow \cdots$$

and

$$\dots \rightarrow H_p(T) \rightarrow H_p(S) \rightarrow H_p(S, T) \rightarrow H_{p-1}(S) \rightarrow H_{p-1}(T) \rightarrow \dots \quad (6.6)$$

The claim follows. □

**Proposition 6.50.** *Let  $T \subset S \subset \mathbb{R}^k$  be a pair of closed and bounded semi-algebraic sets,  $(K, A)$  be a pair of simplicial complexes in  $\mathbb{R}^k$ , with  $A$  a sub-complex of  $K$  and let  $h: |K| \rightarrow S$  be a semi-algebraic homeomorphism such that the image of  $|K|$  is  $T$ . Let  $n_i(K)$  be the number of simplexes of dimension  $i$  of  $K$ , and let  $m_i(A)$  be the number of simplexes of dimension  $i$  of  $A$ . Then*

$$\begin{aligned} \chi(S, T) &= \chi(K, A) \\ &= \sum_i (-1)^i n_i(K) - \sum_i (-1)^i m_i(A). \end{aligned}$$

**Proof:** By Proposition 6.49,  $\chi(K, A) = \chi(K) - \chi(A)$ . The proposition is now a consequence of Proposition 6.3. □

Let  $(X, A), (Y, B)$  be two pairs of semi-algebraic sets. The pairs are  $(X, A), (Y, B)$  are **semi-algebraically homotopy equivalent** if there exist continuous semi-algebraic functions  $f: X \rightarrow Y, g: Y \rightarrow X$  such that,  $\text{Im}(f|_A) \subset B, \text{Im}(g|_B) \subset A$  and such that  $g \circ f \sim \text{Id}_X, g|_B \circ f|_A \sim \text{Id}_A, f \circ g \sim \text{Id}_Y$ , and  $f|_A \circ g|_B \sim \text{Id}_B$ . If two pairs are semi-algebraically homotopy equivalent, we also say that they have the same homotopy type.

We have the following proposition which is a generalization of Proposition 6.35 to pairs of closed, bounded semi-algebraic sets.

**Proposition 6.51.** *Let  $\mathbb{R}$  be a real closed field. Let  $(X, A), (Y, B)$  be two pairs of closed, bounded semi-algebraic sets of  $\mathbb{R}^k$  that are semi-algebraically homotopy equivalent. Then,  $H_*(X, A) \cong H_*(Y, B)$ .*

**Proof:** Since  $(X, A)$  and  $(Y, B)$  are semi-algebraically homotopy equivalent, there exist continuous semi-algebraic functions  $f: X \rightarrow Y, g: Y \rightarrow X$  such that,  $\text{Im}(f|_A) \subset B, \text{Im}(g|_B) \subset A$  and such that  $g \circ f \sim \text{Id}_X, g|_B \circ f|_A \sim \text{Id}_A, f \circ g \sim \text{Id}_Y$ , and  $f|_A \circ g|_B \sim \text{Id}_B$ .

After choosing triangulations of  $X$  and  $Y$  (respecting the subsets  $A$  and  $B$ , respectively) and using the same construction as in the proof of Proposition 6.35, we see that  $f$  induces isomorphisms,  $H_*(f): H_*(X) \rightarrow H_*(Y), H_*(f): H_*(A) \rightarrow H_*(B)$ , as well an homomorphism  $H_*(f): H_*(X, A) \rightarrow H_*(Y, B)$  such that the following diagram commutes.

$$\begin{array}{ccccccccc} H_i(A) & \longrightarrow & H_i(X) & \longrightarrow & H_i(X, A) & \longrightarrow & H_{i-1}(A) & \longrightarrow & H_{i-1}(X) \\ \downarrow H_i(f) & & \downarrow H_i(f) & & \downarrow H_i(f) & & \downarrow H_{i-1}(f) & & \downarrow H_{i-1}(f) \\ H_i(B) & \longrightarrow & H_i(Y) & \longrightarrow & H_i(Y, B) & \longrightarrow & H_{i-1}(B) & \longrightarrow & H_{i-1}(Y) \end{array}$$

The rows correspond to the long exact sequence of the pairs  $(X, A)$  and  $(Y, B)$  (see (6.6)) and the vertical homomorphisms are those induced by  $f$ .

Now applying Lemma 6.11 (Five Lemma) we see that  $H_*(f): H_*(X, A) \rightarrow H_*(Y, B)$  is also an isomorphism.  $\square$

### 6.3.3 Borel-Moore Homology

In this section we will consider **basic locally closed semi-algebraic sets** which are, by definition, intersections of closed semi-algebraic sets with basic open ones. Let  $S \subset \mathbb{R}^k$  be a basic locally closed semi-algebraic set and let  $S_r = S \cap B_k(0, r)$ . The  $p$ -th **Borel-Moore homology group** of  $S$ , denoted by  $H_p^{\text{BM}}(S)$ , is defined to be the  $p$ -th simplicial homology group of the pair  $(\overline{S_r}, \overline{S_r} \setminus S_r)$  for large enough  $r > 0$ . Its dimension is the  $p$ -th Borel-Moore Betti number and is denoted by  $b_p^{\text{BM}}(S)$ . We denote by  $H_*^{\text{BM}}(S)$  the chain complex  $(H_p^{\text{BM}}(S), 0)$  and call it the **Borel-Moore homology of  $S$** .

Note that, for a basic locally closed semi-algebraic set  $S$ , both  $\overline{S_r}$  and  $\overline{S_r} \setminus S_r$  are closed and bounded and hence  $H_i(\overline{S_r}, \overline{S_r} \setminus S_r)$  is well defined. It follows clearly from the definition that for a closed and bounded semi-algebraic set, the Borel-Moore homology groups coincide with the simplicial homology groups.

**Exercise 6.3.** Let  $D$  be the plane minus the origin.

Prove that

$$\begin{aligned} H_0^{\text{BM}}(D) &= \mathbb{Q}, \\ H_1^{\text{BM}}(D) &= \mathbb{Q}, \\ H_2^{\text{BM}}(D) &= \mathbb{Q}. \end{aligned}$$

We will show that the Borel-Moore homology is invariant under semi-algebraic homeomorphisms by proving that the Borel-Moore homology coincides with the simplicial homology of the Alexandrov compactification which we introduce below.

Suppose that  $X \subset \mathbb{R}^k$  is a basic locally closed semi-algebraic set, which is not simultaneously closed and bounded, and that  $X$  is the intersection of a closed semi-algebraic set,  $V$ , with the open semi-algebraic set defined by strict inequalities,  $P_1 > 0, \dots, P_m > 0$ . We will assume that  $X \neq \mathbb{R}^k$ . Otherwise, we embed  $X$  in  $\mathbb{R}^{k+1}$ .

We now define the Alexandrov compactification of  $X$ , denoted by  $\dot{X}$ , having the following properties:

- $\dot{X}$  is closed and bounded,
- there exists a semi-algebraic continuous map,  $\eta: X \rightarrow \dot{X}$ , which is a homeomorphism onto its image,
- $\dot{X} \setminus X$  is a single point.

Let  $T_1, \dots, T_m$  be new variables and  $\pi: \mathbb{R}^{k+m} \rightarrow \mathbb{R}^k$  be the projection map forgetting the new coordinates. Consider the closed semi-algebraic set  $Y \subset \mathbb{R}^{k+m}$  defined as the intersection of  $\pi^{-1}(V)$  with the set defined by

$$T_1^2 P_1 - 1 = \dots = T_m^2 P_m - 1 = 0, \quad T_1 \geq 0, \dots, T_m \geq 0.$$

Clearly,  $Y$  is homeomorphic to  $X$ . After making an affine change of coordinates we can assume that  $Y$  does not contain the origin.

Let  $\phi: \mathbb{R}^{k+m} \setminus \{0\} \rightarrow \mathbb{R}^{k+m}$  be the continuous map defined by

$$\phi(x) = \frac{x}{|x|^2}.$$

We define the **Alexandroff compactification** of  $X$  by

$$\dot{X} = \phi(Y) \cup \{0\},$$

and

$$\eta = \phi \circ \pi|_{\overline{Y}}^{-1}.$$

In case  $X$  is closed and bounded, we define  $\dot{X} = X$ .

We now prove,

**Lemma 6.52.**

- a)  $\eta(X)$  is semi-algebraically homeomorphic to  $X$ ,
- b)  $\dot{X}$  is a closed and bounded semi-algebraic set.

**Proof:** We follow the notations introduced in the definition of  $\dot{X}$ . It is easy to verify that  $\phi$  is a homeomorphism and since  $\pi|_{\overline{Y}}^{-1}$  is also a homeomorphism, it follows that  $\eta$  is a homeomorphism onto its image.

We now prove that  $\dot{X}$  is closed and bounded. It is clear from the definition of  $Y$ , that  $Y$  is a closed and unbounded subset of  $\mathbb{R}^{k+m}$ . Since  $0 \notin Y$ ,  $\phi(Y)$  is bounded. Moreover, if  $x \in \overline{\phi(Y)}$ , but  $x \notin \phi(Y)$ , then  $x = 0$ . Otherwise, if  $x \neq 0$ , then  $\phi^{-1}(x)$  must belong to the closure of  $Y$  and hence to  $Y$  since  $Y$  is closed. But this would imply that  $x \in \phi(Y)$ . This shows that  $\dot{X} = \phi(Y) \cup \{0\}$  is closed.  $\square$

We call  $\dot{X}$  to be the Alexandrov compactification of  $X$ . We now show that the Alexandrov compactification is unique up to semi-algebraic homeomorphisms.

**Theorem 6.53.** *Suppose that  $X$  is as above and  $Z$  is a closed and bounded semi-algebraic set such that,*

- a) *There exists a semi-algebraic continuous map,  $\phi: X \rightarrow Z$ , which gives a homeomorphism between  $X$  and  $\phi(X)$ ,*
- b)  *$Z \setminus \phi(X)$  is a single point.*

Then,  $Z$  is semi-algebraically homeomorphic to  $\dot{X}$ .

**Proof:** Let  $Z = \phi(X) \cup \{z\}$ . We have that  $\dot{X} = \eta(X) \cup \{0\}$ . Since  $\eta(X)$  and  $\phi(X)$  are each homeomorphic to  $X$ , there is an induced homeomorphism,  $\psi: \eta(X) \rightarrow \phi(X)$ . Extend  $\psi$  to  $\dot{X}$  by defining  $\psi(0) = z$ . It is easy to check that this extension is continuous and thus gives a homeomorphism between  $\dot{X}$  and  $Z$ .  $\square$

We finally prove that the Borel-Moore homology groups defined above is invariant under semi-algebraic homeomorphisms.

**Theorem 6.54.** *Let  $X$  be a basic locally closed semi-algebraic set. For any basic locally closed semi-algebraic set  $Y$  which is semi-algebraically homeomorphic to  $X$ , we have that  $H_\star^{\text{BM}}(X) \cong H_\star^{\text{BM}}(Y)$ .*

**Proof:** If  $X$  is closed and bounded there is nothing to prove since,  $\dot{X} = X$  and  $H_\star^{\text{BM}}(X) \cong H_\star(X)$  by definition.

Otherwise, let  $X$  be the intersection of a closed semi-algebraic set,  $V$ , with the open semi-algebraic set defined by strict inequalities,  $P_1 > 0, \dots, P_m > 0$ . We follow the notations used in the definition of the Alexandrov compactification above, as well as those used in the definition of Borel-Moore homology groups.

For  $\varepsilon, \delta > 0$  we define,  $X_{\varepsilon, \delta}$  to be the intersection of  $V \cap B_k(0, \frac{1}{\delta})$  with the set defined by,  $P_1 > \varepsilon, \dots, P_m > \varepsilon$ .

Let  $\dot{X} \subset \mathbb{R}^{k+m}$  be the Alexandrov compactification of  $X$  defined previously, and let  $B_{\varepsilon, \delta} = \overline{B_k(0, \delta)} \times \overline{B_m(0, \varepsilon)} \subset \mathbb{R}^{k+m}$ . It follows from Theorem 5.48 that the pair  $(\dot{X}, 0)$  is homotopy equivalent to  $(\dot{X}, B_{\varepsilon, \delta})$  for all  $0 < \varepsilon \ll \delta \ll 1$ . Moreover, the pair  $(\dot{X}, B_{\varepsilon, \delta})$  is homeomorphic to the pair,  $(\overline{X_{\varepsilon, \delta}}, \overline{X_{\varepsilon, \delta}} \setminus X_{\varepsilon, \delta})$ . It follows again from Theorem 5.48 that the pair,  $(\overline{X_{\varepsilon, \delta}}, \overline{X_{\varepsilon, \delta}} \setminus X_{\varepsilon, \delta})$ , is homotopy equivalent to  $(\overline{X_{0, \delta}}, \overline{X_{0, \delta}} \setminus X_{0, \delta})$ . However, by definition  $H_\star^{\text{BM}}(X) \cong H_\star(\overline{X_{0, \delta}}, \overline{X_{0, \delta}} \setminus X_{0, \delta})$  and hence we have shown that

$$H_\star^{\text{BM}}(X) \cong H_\star(\dot{X}).$$

Since by Theorem 6.53, the Alexandrov compactification is unique up to semi-algebraic homeomorphisms, and by Theorem 6.31 the simplicial homology groups are also invariant under semi-algebraic homeomorphisms, this proves the theorem.  $\square$

We now prove the additivity of Borel-Moore homology groups. More precisely, we prove,

**Theorem 6.55.** *Let  $A \subset X \subset \mathbb{R}^k$  be closed semi-algebraic sets. Then there exists an exact sequence,*

$$\dots \rightarrow H_i^{\text{BM}}(A) \rightarrow H_i^{\text{BM}}(X) \rightarrow H_i^{\text{BM}}(X \setminus A) \rightarrow \dots$$

**Proof:** Let  $Y = X \setminus A$ . By definition,

$H_*^{\text{BM}}(Y) \cong H_*(\overline{Y_r}, \overline{Y_r} \setminus Y_r)$ , where  $Y_r = Y \cap B_k(0, r)$ , with  $r > 0$  and sufficiently large.

Similarly, let  $X_r = X \cap B_k(0, r)$ , and  $A_r = X \cap B_k(0, r)$ . Notice that, since  $X$  and  $A$  are closed,  $\overline{Y_r} \setminus Y_r \subset \overline{A_r}$ . Consider a semi-algebraic triangulation of  $h: |K| \rightarrow \overline{X_r}$  which respects the subset  $\overline{A_r}$ . Let  $K^1 \subset K^2$  denote the sub-complexes corresponding to  $\overline{Y_r} \setminus Y_r \subset \overline{A_r}$ . It is now clear from definition that,  $C_\bullet(K, K^1) \cong C_\bullet(K, K^2)$  and hence  $H_*(K, K^1) \cong H_*(K, K^2)$ . But,

$$H_*(K, K^1) \cong H_*(\overline{Y_r}, \overline{Y_r} \setminus Y_r) \cong H_*^{\text{BM}}(Y)$$

and  $H_*(K, K^2) \cong H_*(\overline{X_r}, \overline{A_r})$ .

This shows that

$$H_*^{\text{BM}}(X \setminus A) \cong H_*(\overline{X_r}, \overline{A_r}).$$

The long exact sequence of homology for the pair  $(\overline{X_r}, \overline{A_r})$  is

$$\cdots \rightarrow H_i(\overline{A_r}) \rightarrow H_i(\overline{X_r}) \rightarrow H_i(\overline{X_r}, \overline{A_r}) \rightarrow \cdots$$

Using the isomorphisms proved above, and the fact that  $X$  (resp.  $A$ ) and  $\overline{X_r}$  (resp.  $\overline{A_r}$ ) are homeomorphic, we get an exact sequence,

$$\cdots \rightarrow H_i(\overline{A}) \rightarrow H_i(\overline{X}) \rightarrow H_i^{\text{BM}}(X \setminus A) \rightarrow \cdots \quad \square$$

### 6.3.4 Euler-Poincaré Characteristic

We define the Euler-Poincaré characteristic for basic locally closed semi-algebraic sets. This definition agrees with the previously defined Euler-Poincaré characteristic for closed and bounded semi-algebraic sets and is additive. The Euler-Poincaré characteristic is a discrete topological invariant of semi-algebraic sets which generalizes the cardinality of a finite set. Hence, additivity is a very natural property to require for Euler-Poincaré characteristic.

We define the **Euler-Poincaré characteristic** of a basic locally closed semi-algebraic set  $S$  by,

$$\chi(S) = \sum_i (-1)^i b_i^{\text{BM}}(S),$$

where  $b_i^{\text{BM}}(S)$  is the dimension of  $H_i^{\text{BM}}(S)$  as a  $\mathbb{Q}$ -vector space. In the special case of a closed and bounded semi-algebraic set, we recover the Euler-Poincaré characteristic already defined.

The Euler-Poincaré characteristic of a basic locally closed semi-algebraic set  $S$ , is related to the Euler-Poincaré characteristic of the closed and bounded semi-algebraic sets  $\overline{S_r}$  and  $\overline{S_r} \setminus S_r$  for all large enough  $r > 0$ , by the following lemma.

**Lemma 6.56.**

$$\chi(S) = \chi(\overline{S_r}) - \chi(\overline{S_r} \setminus S_r),$$

where  $S_r = S \cap B_k(0, r)$  and  $r > 0$  and sufficiently large.

**Proof:** Immediate consequence of the definition and Proposition 6.49 □

**Proposition 6.57. [Additivity of Euler-Poincaré characteristic]** *Let  $S, S_1$  and  $S_2$  be basic locally closed semi-algebraic sets such that  $S_1 \cup S_2 = S$ ,  $S_1 \cap S_2 = \emptyset$ . Then*

$$\chi(S) = \chi(S_1) + \chi(S_2).$$

**Proof:** This is an immediate consequence of Theorem 6.55. □

*Remark 6.58.* Note that the additivity property of the Euler-Poincaré characteristic would not be satisfied if we had defined the Euler-Poincaré characteristic in terms of the homology groups rather than in terms of the Borel-Moore homology groups.

For instance, the Euler-Poincaré of the line would be -1, that of a point would be 1, and that of the line minus the point is 2.

Using the definition of Euler-Poincaré characteristic through Borel-Moore homology, the Euler-Poincaré of the line is 0, that of a point is 1, and that of the line minus the point is -1. □

*Remark 6.59.* Notice that, unlike the ordinary homology (see Proposition 6.47), the Borel-Moore homology is not invariant under semi-algebraic homotopy. For instance, a line is semi-algebraically homotopy equivalent to a point, while their Euler-Poincaré characteristics differ as seen in Remark 6.58. □

Let  $S \subset \mathbb{R}^k$  be a closed semi-algebraic set. Given  $P \in \mathbb{R}[X_1, \dots, X_k]$ , we denote

$$\begin{aligned} \text{Reali}(P = 0, S) &= \{x \in S \mid P(x) = 0\}, \\ \text{Reali}(P > 0, S) &= \{x \in S \mid P(x) > 0\}, \\ \text{Reali}(P < 0, S) &= \{x \in S \mid P(x) < 0\}, \end{aligned}$$

and  $\chi(P = 0, S), \chi(P > 0, S), \chi(P < 0, S)$  the Euler-Poincaré characteristics of the corresponding sets.

The **Euler-Poincaré-query** of  $P$  for  $S$  is

$$\text{EuQ}(P, S) = \chi(P > 0, S) - \chi(P < 0, S).$$

The following equality generalized the basic result of sign determination (Equation 2.6).



**Proposition 6.60.** *The following equality holds:*

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} \chi(P=0, S) \\ \chi(P>0, S) \\ \chi(P<0, S) \end{bmatrix} = \begin{bmatrix} \text{EuQ}(1, S) \\ \text{EuQ}(P, S) \\ \text{EuQ}(P^2, Z) \end{bmatrix} \quad (6.7)$$

**Proof:** We need to prove

$$\begin{aligned} \chi(P=0, S) + \chi(P>0, S) + \chi(P<0, S) &= \text{EuQ}(1, S), \\ \chi(P>0, S) - \chi(P<0, S) &= \text{EuQ}(P, S), \\ \chi(P>0, S) + \chi(P<0, S) &= \text{EuQ}(P^2, S). \end{aligned}$$

The claim is an immediate consequence of Proposition 6.57.  $\square$

## 6.4 Bibliographical Notes

Modern algebraic topology has its origins in the work of Poincaré [130]. The first proof of the independence of the simplicial homology groups from the triangulation of a polyhedron is due to Alexander [2]. The Mayer-Vietoris theorem first occurs in a paper by Vietoris [161]. The Borel-Moore homology groups first appear in [27].

---

## Quantitative Semi-algebraic Geometry

In this chapter, we study various quantitative bounds on the number of connected components and Betti numbers of algebraic and semi-algebraic sets. The key method for this study is the critical point method, i.e. the consideration of the critical points of a well chosen projection. The critical point method also plays a key role for improving the complexity of algorithms in the last chapters of the book.

In Section 7.1, we explain a few basic results of Morse theory and use them to study the topology of a non-singular algebraic hypersurface in terms of the number of critical points of a well chosen projection. Bounding the number of these critical points by Bézout's theorem provides a bound on the sum of the Betti numbers of a non-singular bounded algebraic hypersurface in Section 7.2. Then we prove a similar bound on the sum of the Betti numbers of a general algebraic set.

In Section 7.3, we prove a bound on the sum of the  $i$ -th Betti numbers over all realizable sign conditions of a finite set of polynomials. In particular, the bound on the zero-th Betti numbers gives us a bound on the number of realizable sign conditions of a finite set of polynomials. We also explain why these bounds are reasonably tight.

In Section 7.4, we prove bounds on Betti numbers of closed semi-algebraic sets. In Section 7.5 we prove that any semi-algebraic set is semi-algebraically homotopic to a closed one and prove bounds on Betti numbers of general semi-algebraic sets.

### 7.1 Morse Theory

We first define the kind of hypersurfaces we are going to consider.

A **non-singular algebraic hypersurface** is the zero set  $\text{Zer}(Q, \mathbb{R}^k)$  of a polynomial  $Q \in \mathbb{R}[X_1, \dots, X_k]$  such that the **gradient** of  $Q$ , i.e. the vector

$$\text{Grad}(Q)(p) = \left( \frac{\partial Q}{\partial X_1}(p), \dots, \frac{\partial Q}{\partial X_k}(p) \right) \text{ is never } 0 \text{ for } p \in \text{Zer}(Q, \mathbb{R}^k).$$

**Exercise 7.1.** Prove that a non-singular algebraic hypersurface is an  $\mathcal{S}^\infty$  submanifold of dimension  $k - 1$ . (Hint. Use the Semi-algebraic implicit function theorem (Theorem 3.25).)

**Exercise 7.2.** Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular algebraic hypersurface. Prove that the gradient vector of  $Q$  at a point  $p \in \text{Zer}(Q, \mathbb{R}^k)$  is orthogonal to the tangent space  $T_p(\text{Zer}(Q, \mathbb{R}^k))$  to  $\text{Zer}(Q, \mathbb{R}^k)$  at  $p$ .

We denote by  $\pi$  the projection from  $\mathbb{R}^k$  to the first coordinate sending  $(x_1, \dots, x_k)$  to  $x_1$ .

**Notation 7.1. [Fiber]** For  $S \subset \mathbb{R}^k$ ,  $X \subset \mathbb{R}$ , let  $S_X$  denote  $S \cap \pi^{-1}(X)$ . We also use the abbreviations  $S_x$ ,  $S_{<x}$ , and  $S_{\leq x}$  for  $S_{\{x\}}$ ,  $S_{(-\infty, x)}$ , and  $S_{(-\infty, x]}$ .  $\square$

Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular algebraic hypersurface and  $p \in \text{Zer}(Q, \mathbb{R}^k)$ . Then, the derivative  $d\pi(p)$  of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  is a linear map from  $T_p(\text{Zer}(Q, \mathbb{R}^k))$  to  $\mathbb{R}$ . Clearly,  $p$  is a critical point of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  if and only if

$$\frac{\partial Q}{\partial X_i}(p) = 0, 2 \leq i \leq k$$

(see Definition 5.55). In other words,  $p$  is a critical point of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  if and only if the gradient of  $Q$  is parallel to the  $X_1$ -axis, i.e.  $T_p(\text{Zer}(Q, \mathbb{R}^k))$  is orthogonal to the  $X_1$  direction. A critical value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  is the projection to the  $X_1$ -axis of a critical point of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ .

**Lemma 7.2.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface. The set of values that are not critical for  $\pi$  is non-empty and open.*

**Proof:** The set of values that are not critical for  $\pi$  is clearly open, from the definition of a critical value. It is also non-empty by Theorem 5.56 (Sard's theorem) since the set of critical values is a finite subset of  $\mathbb{R}$ .  $\square$

Also, as an immediate consequence of the Semi-algebraic implicit function theorem (Theorem 3.25), we have:

**Proposition 7.3.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface. If  $x$  is not a critical value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  and  $p$  is a point of  $\text{Zer}(Q, \mathbb{R}^k)_x$ , then for  $\epsilon$  small enough  $\text{Zer}(Q, \mathbb{R}^k) \cap B(p, \epsilon)_{<x}$  is non-empty and semi-algebraically connected.*

We also have the following proposition.

**Proposition 7.4.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface. The set of critical points of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  meets every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof:** Let  $C$  be a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ . Since  $C$  is semi-algebraic, closed, and bounded, its image by  $\pi$  is semi-algebraic, closed, and bounded, using Theorem 3.20. Thus  $\pi(C)$  is a finite number of points and closed intervals and has a smallest element  $v$ . Using Proposition 7.3, it is clear that any  $x \in C$  such that  $\pi(x) = v$  is critical.  $\square$

We will now state and prove the first basic ingredient of Morse theory. In the remainder of the section, we assume  $\mathbb{R} = \mathbb{R}$ . We suppose that  $\text{Zer}(Q, \mathbb{R}^k)$  is a bounded algebraic non-singular hypersurface and denote by  $\pi$  the projection map sending  $(x_1, \dots, x_k)$  to  $x_1$ .

Consider the sets  $\text{Zer}(Q, \mathbb{R}^k)_{\leq x}$  as  $x$  varies from  $-\infty$  to  $\infty$ . Thinking of  $X_1$  as the horizontal axis, the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq x}$  is the part of  $\text{Zer}(Q, \mathbb{R}^k)$  to the left of the hyperplane defined by  $X_1 = x$ , and we study the changes in the homotopy type of this set as we sweep the hyperplane in the rightward direction. Theorem 7.5 states that there is no change in the homotopy type as  $x$  varies strictly between two critical values of  $\pi$ .

**Theorem 7.5. [Morse lemma A]** *Let  $[a, b]$  be an interval containing no critical value of  $\pi$ . Then  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  and  $\text{Zer}(Q, \mathbb{R}^k)_a \times [a, b]$  are homeomorphic, and  $\text{Zer}(Q, \mathbb{R}^k)_{\leq a}$  is homotopy equivalent to  $\text{Zer}(Q, \mathbb{R}^k)_{\leq b}$ .*

Theorem 7.5 immediately implies:

**Proposition 7.6.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface,  $[a, b]$  such that  $\pi$  has no critical value in  $[a, b]$ , and  $d \in [a, b]$ .*

- *The sets  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  and  $\text{Zer}(Q, \mathbb{R}^k)_d$  have the same number of semi-algebraically connected components.*
- *Let  $S$  be a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . Then, for every  $d \in [a, b]$ ,  $S_d$  is semi-algebraically connected.*

The proof of Theorem 7.5 is based on local existence and uniqueness of solutions to systems of differential equations. Let  $U$  be an open subset of  $\mathbb{R}^k$ . A **vector field**  $\Gamma$  on  $U$  is a  $C^\infty$  map from an open set  $U$  of  $\mathbb{R}^k$  to  $\mathbb{R}^k$ . To a vector field is associated a system of differential equations

$$\frac{dx_i}{dT} = \Gamma_i(x_1, \dots, x_k), 1 \leq i \leq k.$$

A **flow line** of the vector field  $\Gamma$  is a  $C^\infty$  map  $\gamma: I \rightarrow \mathbb{R}^k$  defined on some interval  $I$  and satisfying

$$\frac{d\gamma}{dT}(t) = \Gamma(\gamma(t)), t \in I.$$

**Theorem 7.7.** *Let  $\Gamma$  be a vector field on an open subset  $V$  of  $\mathbb{R}^k$  such that for every  $x \in V$ ,  $\Gamma(x) \neq 0$ . For every  $y \in V$ , there exists a neighborhood  $U$  of  $y$  and  $\epsilon > 0$ , such that for every  $x \in U$ , there exists a unique flow line  $\gamma_x: (-\epsilon, \epsilon) \rightarrow \mathbb{R}^k$  of  $\Gamma$  satisfying the initial condition  $\gamma_x(0) = x$ .*

**Proof:** Since  $\Gamma$  is  $C^\infty$ , there exists a bounded neighborhood  $W$  of  $y$  and  $L > 0$  such that  $|\Gamma(x_1) - \Gamma(x_2)| < L|x_1 - x_2|$  for all  $x_1, x_2 \in W$ . Let  $A = \sup_{x \in W} |\Gamma(x)|$ . Also, let  $\epsilon > 0$  be a small enough number such that the set

$$W' = \{x \in W \mid B_k(x, \epsilon A) \subset W\}$$

contains an open set  $U$  containing  $y$ .

Let  $x \in U$ . If  $\gamma_x: [-\epsilon, \epsilon] \rightarrow \mathbb{R}^k$ , with  $\gamma_x(0) = x$ , is a solution, then  $\gamma_x([-\epsilon, \epsilon]) \subset W'$ . This is because  $|\Gamma(x')| \leq A$  for every  $x' \in W$ , and hence applying the Mean Value Theorem,  $|x - \gamma_x(t)| \leq |t|A$  for all  $t \in [-\epsilon, \epsilon]$ . Now, since  $x \in U$ , it follows that  $\gamma_x([-\epsilon, \epsilon]) \subset W'$ .

We construct the solution  $\gamma_x: [-\epsilon, \epsilon] \rightarrow W$  as follows. Let  $\gamma_{x,0}(t) = x$  for all  $t$  and

$$\gamma_{x,n+1}(t) = x + \int_0^t \Gamma(\gamma_{x,n}(t)) dt.$$

Note that  $\gamma_{x,n}([-\epsilon, \epsilon]) \subset W'$  for every  $n \geq 0$ . Now,

$$\begin{aligned} |\gamma_{x,n+1}(t) - \gamma_{x,n}(t)| &= \left| \int_0^t (\Gamma(\gamma_{x,n}(t)) - \Gamma(\gamma_{x,n-1}(t))) dt \right| \\ &\leq \left| \int_0^t |\Gamma(\gamma_{x,n}(t)) - \Gamma(\gamma_{x,n-1}(t))| dt \right| \\ &\leq \epsilon L |\gamma_{x,n}(t) - \gamma_{x,n-1}(t)| \end{aligned}$$

Choosing  $\epsilon$  such that  $\epsilon < 1/L$ , we see that for every fixed  $t \in [-\epsilon, \epsilon]$ , the sequence  $\gamma_{x,n}(t)$  is a Cauchy sequence and converges to a limit  $\gamma_x(t)$ .

Moreover, it is easy to verify that  $\gamma_x(t)$  satisfies the equation,

$$\gamma_x(t) = x + \int_0^t \Gamma(\gamma_x(t)) dt.$$

Differentiating both sides, we see that  $\gamma_x(t)$  is a flow line of the given vector field  $\Gamma$ , and clearly  $\gamma_x(0) = x$ .

The proof of uniqueness is left as an exercise.  $\square$

Given a  $C^\infty$  hypersurface  $M \subset \mathbb{R}^k$ , a  $C^\infty$  **vector field on  $M$** ,  $\Gamma$ , is a  $C^\infty$  map that associates to each  $x \in M$  a tangent vector  $\Gamma(x) \in T_x(M)$ .

An important example of a vector field on a hypersurface is the gradient vector field. Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular algebraic hypersurface and  $(a', b')$  such that  $\pi$  has no critical point on  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$ . The **gradient vector field of  $\pi$**  on  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$  is the  $C^\infty$  vector field on  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$  that to every  $p \in \text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$  associates  $\Gamma(p)$  characterized by the following properties

- it belongs to  $T_p(\text{Zer}(Q, \mathbb{R}^k))$ ,

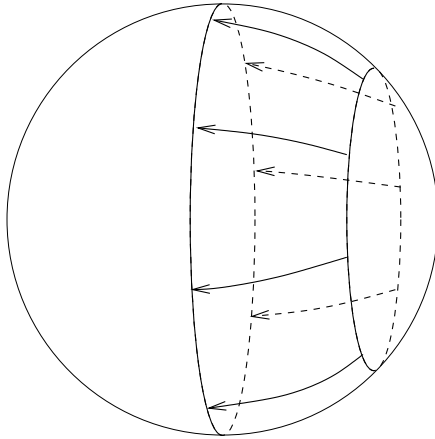
- it belongs to the plane generated by the gradient  $\text{Grad}(Q)(p)$ , and the unit vector of the  $X_1$ -axis,
- its projection on the  $X_1$ -axis is the negative of the unit vector.

The flow lines of the gradient vector field correspond to curves on the hypersurface along which the  $X_1$  coordinate decreases maximally. A straightforward computation shows that, for  $p \in \text{Zer}(Q, \mathbb{R}^k)$ ,

$$\Gamma(p) = - \frac{G(p)}{\sum_{2 \leq i \leq k} \left( \frac{\partial Q}{\partial X_i}(p) \right)^2},$$

where

$$G(p) = \left( \sum_{2 \leq i \leq k} \left( \frac{\partial Q}{\partial X_i}(p) \right)^2, -\frac{\partial Q}{\partial X_1} \frac{\partial Q}{\partial X_2}(p), \dots, -\frac{\partial Q}{\partial X_1} \frac{\partial Q}{\partial X_k}(p) \right).$$



**Fig. 7.1.** Flow of the gradient vector field on the 2-sphere

**Proof of Theorem 7.5:** By Lemma 7.2 we can chose  $a' < a, b' > b$  such that  $\pi$  has no critical point on  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$ . Consider the gradient vector field of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$ .

By Corollary 5.51, the set  $\text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$  can be covered by a finite number of open sets such that for each open set  $U'$  in the cover, there is an open  $U$  of  $\mathbb{R}^{k-1}$  and a diffeomorphism  $\Phi: U \rightarrow U'$ .

Using the linear maps  $d\Phi_x^{-1}: T_x(M) \rightarrow T_{\Phi^{-1}(x)}\mathbb{R}^{k-1}$ , we associate to the gradient vector field of  $\pi$  on  $U' \subset \text{Zer}(Q, \mathbb{R}^k)_{(a', b')}$  a  $C^\infty$  vector field on  $U$ .

By Theorem 7.7, for each point  $x \in \text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ , there exists a neighborhood  $W$  of  $\Phi^{-1}(x)$  and an  $\epsilon > 0$  such that the induced vector field in  $\mathbb{R}^{k-1}$  has a solution  $\gamma_x(t)$  for  $t \in (-\epsilon, \epsilon)$  and such that  $\gamma_x(0) = \Phi^{-1}(x)$ . We consider its image  $\Phi(\gamma_x)$  on  $\text{Zer}(Q, \mathbb{R}^k)_{(a',b')}$ . Thus, for each point  $x \in \text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ , we have a neighborhood  $U_x$ , a number  $\epsilon_x > 0$ , and a curve

$$\Phi \circ \gamma_x: (-\epsilon_x, \epsilon_x) \rightarrow \text{Zer}(Q, \mathbb{R}^k)_{(a',b')},$$

such that  $\Phi \circ \gamma_x(0) = x$ , and  $d\Phi \circ \gamma_x dt = \Gamma(\Phi \circ \gamma_x(t))$  for all  $t \in (-\epsilon_x, \epsilon_x)$ .

Since  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  is compact, we can cover it using a finite number of the neighborhoods  $U_x$  and let  $\epsilon_0 > 0$  be the least among the corresponding  $\epsilon_x$ 's. For  $t \in [0, b - a]$ , we define a one-parameter family of smooth maps

$$\alpha_t: \text{Zer}(Q, \mathbb{R}^k)_b \rightarrow \text{Zer}(Q, \mathbb{R}^k)_{\leq b}$$

as follows:

Let  $x \in \text{Zer}(Q, \mathbb{R}^k)_b$ . If  $|t| \leq \epsilon_0/2$ , we let  $\alpha_t(x) = \gamma_x(t)$ . If  $|t| > \epsilon_0/2$ , we write  $t = n \epsilon_0/2 + \delta$ , where  $n$  is an integer and  $|\delta| < \epsilon_0/2$ . We let

$$\alpha_t(x) = \overbrace{\alpha_{\epsilon_0/2} \circ \dots \circ \alpha_{\epsilon_0/2}}^{n \text{ times}} \circ \alpha_\delta(x).$$

Observe the following.

- For every  $x \in \text{Zer}(Q, \mathbb{R}^k)_b$ ,  $\alpha_0(x) = x$ .
- By construction,  $\frac{d\alpha_t(x)}{dt} = \Gamma(\alpha_t(x))$ . Since the projection on the  $X_1$  axis of  $\Gamma(\alpha_t(x)) = (-1, 0, \dots, 0)$ , it follows that  $\pi(\alpha_t(x)) = b - t$ .
- $\alpha_t(\text{Zer}(Q, \mathbb{R}^k)_b) = \text{Zer}(Q, \mathbb{R}^k)_{b-t}$ .
- It follows from the uniqueness of the flowlines through every point of the gradient vector field on  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  (Theorem 7.7) that each  $\alpha_t$  defined above is injective.

We now claim that the map  $f: \text{Zer}(Q, \mathbb{R}^k)_{[a,b]} \rightarrow \text{Zer}(Q, \mathbb{R}^k)_a \times [a, b]$  defined by

$$f(x) = (\alpha_{b-a}(\alpha_{b-\pi(x)}^{-1}(x)), \pi(x))$$

is a homeomorphism. This is an immediate consequence of the properties of  $\alpha_t$  listed above.

Next, consider the map  $F(x, t): \text{Zer}(Q, \mathbb{R}^k)_{\leq b} \times [0, 1] \rightarrow \text{Zer}(Q, \mathbb{R}^k)_{\leq b}$  defined as follows:

$$\begin{aligned} F(x, s) &= x, && \text{if } \pi(x) \leq b - s(b - a) \\ &= \alpha_{s(b-a)}(\alpha_{b-\pi(x)}^{-1}(x)), && \text{otherwise.} \end{aligned}$$

Clearly,  $F$  is a deformation retraction from  $\text{Zer}(Q, \mathbb{R}^k)_{\leq b}$  to  $\text{Zer}(Q, \mathbb{R}^k)_{\leq a}$ , so that  $\text{Zer}(Q, \mathbb{R}^k)_{\leq b}$  is homotopy equivalent to  $\text{Zer}(Q, \mathbb{R}^k)_{\leq a}$ .

This completes the proof. □

Theorem 7.5 states that there is no change in homotopy type on intervals containing no critical values. The remainder of the section is devoted to studying the changes in homotopy type that occur at the critical values. In this case, we will not be able to use the gradient vector field of  $\pi$  to get a flow as the gradient becomes zero at a critical point. We will, however, show how to modify the gradient vector field in a neighborhood of a critical point so as to get a new vector field that agrees with the gradient vector field outside a small neighborhood. The flow corresponding to this new vector field will give us a homotopy equivalence between  $\text{Zer}(Q, \mathbb{R}^k)_{\leq c+\epsilon}$  and  $\text{Zer}(Q, \mathbb{R}^k)_{\leq c-\epsilon} \cup B$ , where  $c$  is a critical value of  $\pi$ ,  $\epsilon > 0$  is sufficiently small, and  $B$  a topological ball attached to  $\text{Zer}(Q, \mathbb{R}^k)_{\leq c-\epsilon}$  by its boundary. The key notion necessary to work this idea out is that of a Morse function.

**Definition 7.8. [Morse function]** Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface and  $\pi$  the projection on the  $X_1$ -axis sending  $x = (x_1, \dots, x_k) \in \mathbb{R}^k$  to  $x_1 \in \mathbb{R}$ . Let  $p \in \text{Zer}(Q, \mathbb{R}^k)$  be a critical point of  $\pi$ . The tangent space  $T_p(\text{Zer}(Q, \mathbb{R}^k))$  is the  $(k - 1)$ -dimensional space spanned by the  $X_2, \dots, X_k$  coordinates with origin  $p$ . By virtue of the Implicit Function Theorem (Theorem 3.25), we can choose  $(X_2, \dots, X_k)$  to be a local system of coordinates in a sufficiently small neighborhood of  $p$ . More precisely, we have an open neighborhood  $U \subset \mathbb{R}^{k-1}$  of  $p' = (p_2, \dots, p_k)$  and a mapping  $\phi: U \rightarrow \mathbb{R}$ , such that, with  $x' = (x_2, \dots, x_k)$ , and

$$\Phi(x') = (\phi(x'), x') \in \text{Zer}(Q, \mathbb{R}^k), \quad (7.1)$$

the mapping  $\Phi$  is a diffeomorphism from  $U$  to  $\Phi(U)$ .

The critical point  $p$  is **non-degenerate** if the  $(k - 1) \times (k - 1)$  Hessian matrix

$$\text{Hes}_\pi(p') = \left[ \frac{\partial^2 \phi}{\partial X_i \partial X_j}(p') \right], \quad 2 \leq i, j \leq k, \quad (7.2)$$

is invertible. Note that  $\text{Hes}_\pi(p')$  is a real symmetric matrix and hence all its eigenvalues are real (Theorem 4.42). Moreover, if  $p$  is a non-degenerate critical point, then all eigenvalues are non-zero. The number of positive eigenvalues of  $\text{Hes}_\pi(p')$  is the **index** of the critical point  $p$ .

The function  $\pi$  is a **Morse function** if all its critical points are non-degenerate and there is at most one critical point of  $\pi$  above each  $x \in \mathbb{R}$ .  $\square$

We next show that to require  $\pi$  to be a Morse function is not a big loss of generality, since an orthogonal change of coordinates can make the projection map  $\pi$  a Morse function on  $\text{Zer}(Q, \mathbb{R}^k)$ .

**Proposition 7.9.** *Up to an orthogonal change of coordinates, the projection  $\pi$  to the  $X_1$ -axis is a Morse function.*

The proof of Proposition 7.9 requires some preliminary work.

We start by proving:



**Proposition 7.10.** *Let  $d$  be the degree of  $Q$ . Suppose that the projection  $\pi$  on the  $X_1$ -axis has only non-degenerate critical points. The number of critical points of  $\pi$  is finite and bounded by  $d(d - 1)^{k-1}$ .*

**Proof:** The critical points of  $\pi$  can be characterized as the real solutions of the system of  $k$  polynomial equations in  $k$  variables

$$Q = \frac{\partial Q}{\partial X_2} = 0, \dots, \frac{\partial Q}{\partial X_k} = 0.$$

We claim that every real solution  $p$  of this system is non-singular, i.e. the Jacobian matrix

$$\begin{bmatrix} \frac{\partial Q}{\partial X_1}(p) & \frac{\partial^2 Q}{\partial X_2 \partial X_1}(p) & \cdots & \frac{\partial^2 Q}{\partial X_k \partial X_1}(p) \\ \vdots & \vdots & & \vdots \\ \frac{\partial Q}{\partial X_k}(p) & \frac{\partial^2 Q}{\partial X_2 \partial X_k}(p) & \cdots & \frac{\partial^2 Q}{\partial X_k \partial X_k}(p) \end{bmatrix}$$

is non-singular. Differentiating the identity (7.3) and evaluating at  $p$ , we obtain for  $2 \leq i, j \leq k$ , with  $p' = (p_2, \dots, p_k)$ ,

$$\frac{\partial^2 Q}{\partial X_j \partial X_i}(p) = -\frac{\partial Q}{\partial X_1}(p) \frac{\partial^2 \phi}{\partial X_j \partial X_i}(p').$$

Since  $\frac{\partial Q}{\partial X_1}(p) \neq 0$  and  $\frac{\partial Q}{\partial X_i}(p) = 0$ , for  $2 \leq i \leq k$ , the claim follows. By Theorem 4.106 (Bézout’s theorem), the number of critical points of  $\pi$  is less than or equal to the product

$$\deg(Q) \deg\left(\frac{\partial Q}{\partial X_2}\right) \cdots \deg\left(\frac{\partial Q}{\partial X_k}\right) = d(d - 1)^{k-1}. \quad \square$$

We interpret geometrically the notion of non-degenerate critical point.

**Proposition 7.11.** *Let  $p \in \text{Zer}(Q, \mathbb{R}^k)$  be a critical point of  $\pi$ . Let  $g: \text{Zer}(Q, \mathbb{R}^k) \rightarrow S^{k-1}(0, 1)$  be the Gauss map defined by*

$$g(x) = \frac{\text{Grad}(Q(x))}{\|\text{Grad}(Q(x))\|}.$$

*The Gauss map is an  $\mathcal{S}^\infty$ -diffeomorphism in a neighborhood of  $p$  if and only if  $p$  is a non-degenerate critical point.*

**Proof:** Since  $p$  is a critical point of  $\pi$ ,  $g(p) = (\pm 1, 0, \dots, 0)$ . Using Notation 7.8, for  $x' \in U$ ,  $x = \Phi(x') = (\phi(x'), x')$ , and applying the chain rule,

$$\frac{\partial Q}{\partial X_i}(x) + \frac{\partial Q}{\partial X_1}(x) \frac{\partial \phi}{\partial X_i}(x') = 0, \quad 2 \leq i \leq k. \tag{7.3}$$

Thus

$$g(x) = \pm \frac{1}{\sqrt{1 + \sum_{i=2}^k \left( \frac{\partial \phi}{\partial X_i}(x') \right)^2}} \left( -1, \frac{\partial \phi}{\partial X_2}(x'), \dots, \frac{\partial \phi}{\partial X_k}(x') \right).$$

Taking the partial derivative with respect to  $X_i$  of the  $j$ -th coordinate  $g_j$  of  $g$ , for  $2 \leq i, j \leq k$ , and evaluating at  $p$ , we obtain

$$\frac{\partial g_j}{\partial X_i}(p) = \pm \frac{\partial^2 \phi}{\partial X_j \partial X_i}(p'), \quad 2 \leq i, j \leq k.$$

The matrix  $[\partial g_i / \partial X_j(p)]$ ,  $2 \leq i, j \leq k$ , is invertible if and only if  $p$  is a non-degenerate critical point of  $\phi$  by (7.2).  $\square$

**Proposition 7.12.** *Up to an orthogonal change of coordinates, the projection  $\pi$  to the  $X_1$ -axis has only non-degenerate critical points.*

**Proof:** Consider again the Gauss map  $g: \text{Zer}(Q, \mathbb{R}^k) \rightarrow S^{k-1}(0, 1)$ , defined by

$$g(x) = \frac{\text{Grad}(Q(x))}{\|\text{Grad}(Q(x))\|}.$$

According to Sard's theorem (Theorem 5.56) the dimension of the set of critical values of  $g$  is at most  $k - 2$ . We prove now that there are two antipodal points of  $S^{k-1}(0, 1)$  such that neither is a critical value of  $g$ . Assume the contrary and argue by contradiction. Since the dimension of the set of critical values is at most  $k - 2$ , there exists a non-empty open set  $U$  of regular values in  $S^{k-1}(0, 1)$ . The set of points that are antipodes to points in  $U$  is non-empty, open in  $S^{k-1}(0, 1)$  and all critical, contradicting the fact that the critical set has dimension at most  $k - 2$ .

After rotating the coordinate system, we may assume that  $(1, 0, \dots, 0)$  and  $(-1, 0, \dots, 0)$  are not critical values of  $g$ . The claim follows from Proposition 7.11.  $\square$

It remains to prove that it is possible to ensure, changing the coordinates if necessary, that there is at most one critical point of  $\pi$  above each  $x \in \mathbb{R}$ .

Suppose that the projection  $\pi$  on the  $X_1$ -axis has only non-degenerate critical points. These critical points are finite in number according to Proposition 7.10. We can suppose without loss of generality that all the critical points have distinct  $X_2$  coordinates, making if necessary an orthogonal change of coordinates in the variables  $X_2, \dots, X_k$  only.

**Lemma 7.13.** *Let  $\delta$  be a new variable and consider the field  $\mathbb{R}\langle\delta\rangle$  of algebraic Puiseux series in  $\delta$ . The set  $S$  of points  $\bar{p} = (\bar{p}_1, \dots, \bar{p}_k) \in \text{Zer}(Q, \mathbb{R}\langle\delta\rangle^k)$  with gradient vector  $\text{Grad}(Q)(\bar{p})$  proportional to  $(1, \delta, 0, \dots, 0)$  is finite. Its number of elements is equal to the number of critical points of  $\pi$ . Moreover there is a point  $\bar{p}$  of  $S$  infinitesimally close to every critical point  $p$  of  $\pi$  and the signature of the Hessian at  $p$  and  $\bar{p}$  coincide.*

**Proof:** Note that, modulo the orthogonal change of variable

$$X'_1 = X_1 + \delta X_2, X'_2 = X_2 - \delta X_1, X'_i = X_i, i \geq 3,$$

a point  $\bar{p}$  such that  $\text{Grad}(Q)(\bar{p})$  is proportional to  $(1, \delta, 0, \dots, 0)$  is a critical point of the projection  $\pi'$  on the  $X'_1$ -axis, and the corresponding critical value of  $\pi'$  is  $\bar{p}_1 + \delta \bar{p}_2$ .

Since  $\text{Zer}(Q, \mathbb{R}^k)$  is bounded, a point  $\bar{p} \in \text{Zer}(Q, \mathbb{R}\langle\delta\rangle^k)$  always has an image by  $\lim_\delta$ . If  $\bar{p}$  is such that  $\text{Grad}(Q)(\bar{p})$  is proportional to  $(1, \delta, 0, \dots, 0)$ , then  $\text{Grad}(Q)(\lim_\delta(\bar{p}))$  is proportional to  $(1, 0, \dots, 0, 0)$ , and thus  $p = \lim_\delta(\bar{p})$  is a critical point of  $\pi$ . Suppose without loss of generality that  $\text{Grad}(Q)(p) = (1, 0, \dots, 0, 0)$ . Since  $p$  is a non-degenerate critical point of  $\pi$ , Proposition 7.11 implies that there is a semi-algebraic neighborhood  $U$  of  $p' = (p_2, \dots, p_k)$  such that  $g \circ \Phi$  is a diffeomorphism from  $U$  to a semi-algebraic neighborhood of  $(1, 0, \dots, 0, 0) \in S^{k-1}(0, 1)$ . Denoting by  $g'$  the inverse of the restriction of  $g$  to  $\Phi(U)$  and considering

$$\text{Ext}(g', \mathbb{R}\langle\delta\rangle) : \text{Ext}(g(\Phi(U)), \mathbb{R}\langle\delta\rangle) \rightarrow \text{Ext}(\Phi(U), \mathbb{R}\langle\delta\rangle),$$

there is a unique  $\bar{p} \in \text{Ext}(\Phi(U), \mathbb{R}\langle\delta\rangle)$  such that  $\text{Grad}(Q)(\bar{p})$  is proportional to  $(1, \delta, 0, \dots, 0)$ . Moreover, denoting by  $J$  the Jacobian of  $\text{Ext}(g', \mathbb{R}\langle\delta\rangle)$ , the value  $J(1, 0, 0, \dots, 0) = t$  is a non-zero real number. Thus the signature of the Hessian at  $p$  and  $\bar{p}$  coincide.  $\square$

**Proof of Proposition 7.9:** Since  $J$  is the Jacobian of  $\text{Ext}(g', \mathbb{R}\langle\delta\rangle)$ , the value  $J(1, 0, 0, \dots, 0) = t$  is a non-zero real number,  $\lim_\delta(J(y)) = t$  for every  $y \in \text{Ext}(S^{k-1}(0, 1), \mathbb{R}\langle\delta\rangle)$  infinitesimally close to  $(1, 0, 0, \dots, 0)$ . Using the mean value theorem (Corollary 2.23)

$$o(|\bar{p} - p|) = o\left(\left|\frac{1}{\sqrt{1 + \delta^2}}(1, \delta, 0, \dots, 0) - (1, 0, 0, \dots, 0)\right|\right) = 1.$$

Thus  $o(\bar{p}_i - p_i) \geq 1, i \geq 1$ .

Let  $b_{i,j} = \frac{\partial^2 \phi}{\partial X_i \partial X_j}(p), 2 \leq i \leq k, 2 \leq j \leq k$ . Taylor's formula at  $p$  for  $\phi$  gives

$$\bar{p}_1 = p_1 + \sum_{2 \leq i \leq k, 2 \leq j \leq k} b_{i,j} (\bar{p}_i - p_i) (\bar{p}_j - p_j) + c,$$

with  $o(c) \geq 2$ . Thus  $o(\bar{p}_1 - p_1) \geq 2$ .

It follows that the critical value of  $\pi'$  at  $\bar{p}$  is  $\bar{p}_1 + \delta \bar{p}_2 = p_1 + \delta p_2 + w$ , with  $o(w) > 1$ .

Thus, all the critical values of  $\pi'$  on  $\text{Zer}(Q, \mathbb{R}\langle\delta\rangle^k)$  are distinct since all values of  $p_2$  are. Using Proposition 3.17, we can replace  $\delta$  by  $d \in \mathbb{R}$ , and we have proved that there exists an orthogonal change of variable such that  $\pi$  is a Morse function.  $\square$

We are now ready to state the second basic ingredient of Morse theory, which is describing precisely the change in the homotopy type that occurs in  $\text{Zer}(Q, \mathbb{R}^k)_{\leq x}$  as  $x$  crosses a critical value when  $\pi$  is a Morse function.

**Theorem 7.14. [Morse lemma B]** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface such that the projection  $\pi$  to the  $X_1$ -axis is a Morse function. Let  $p$  be a non-degenerate critical point of  $\pi$  of index  $\lambda$  and such that  $\pi(p) = c$ .*

*Then, for all sufficiently small  $\epsilon > 0$ , the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq c+\epsilon}$  has the homotopy type of the union of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq c-\epsilon}$  with a ball of dimension  $k - 1 - \lambda$ , attached along its boundary.*

We first prove a lemma that will allow us to restrict to the case where  $x = 0$  and where  $Q$  is a quadratic polynomial of a very simple form.

Let  $\text{Zer}(Q, \mathbb{R}^k), U, \phi, \Phi$  be as above (see page 243).

**Lemma 7.15.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface such that the projection  $\pi$  to the  $X_1$ -axis is a Morse function. Let  $p \in \text{Zer}(Q, \mathbb{R}^k)$  be a non-degenerate critical point of the map  $\pi$  with index  $\lambda$ . Then there exists an open neighborhood  $V$  of the origin in  $\mathbb{R}^{k-1}$  and a diffeomorphism  $\Psi$  from  $U$  to  $V$  such that, denoting by  $Y_i$  the  $i$ -th coordinate of  $\Psi(X_2, \dots, X_k)$ ,*

$$\phi(Y_2, \dots, Y_k) = \sum_{2 \leq i \leq \lambda+1} Y_i^2 - \sum_{\lambda+2 \leq i \leq k} Y_i^2.$$

**Proof:** We assume without loss of generality that  $p$  is the origin. Also, by Theorem 4.42, we assume that the matrix

$$\text{Hes}(0) = \left[ \frac{\partial^2 \phi}{\partial X_i \partial X_j}(0) \right], \quad 2 \leq i, j \leq k,$$

is diagonal with its first  $\lambda$  entries  $+1$  and the remaining  $-1$ .

Let us prove that there exists a  $C^\infty$  map  $M$  from  $U$  to the space of symmetric  $(k - 1) \times (k - 1)$  matrices,  $X \mapsto M(X) = (m_{ij}(X))$ , such that

$$\phi(X_2, \dots, X_k) = \sum_{2 \leq i, j \leq k} m_{ij}(X) X_i X_j.$$

Using the fundamental theorem of calculus twice, we obtain

$$\begin{aligned} \phi(X_2, \dots, X_k) &= \sum_{2 \leq j \leq k} X_j \int_0^1 \frac{\partial \phi}{\partial X_j}(t X_2, \dots, t X_k) dt \\ &= \sum_{2 \leq i \leq k} \sum_{2 \leq j \leq k} X_i X_j \int_0^1 \int_0^1 \frac{\partial^2 \phi}{\partial X_i \partial X_j}(s t X_2, \dots, s t X_k) dt ds. \end{aligned}$$

Take

$$m_{ij}(X_2, \dots, X_k) = \int_0^1 \int_0^1 \frac{\partial^2 \phi}{\partial X_i \partial X_j}(s t X_2, \dots, s t X_k) dt ds.$$

Note that the matrix  $M(X_2, \dots, X_k)$  obtained above clearly satisfies  $M(0) = H(0)$ , and  $M(x_2, \dots, x_k)$  is close to  $H(x_2, \dots, x_k)$  for  $(x_2, \dots, x_k)$  in a sufficiently small neighborhood of the origin.

Using Theorem 4.42 again, there exists a  $C^\infty$  map  $N$  from a sufficiently small neighborhood  $V$  of 0 in  $\mathbb{R}^{k-1}$  to the space of  $(k-1) \times (k-1)$  real invertible matrices such that

$$\forall x \in V, N(x)^t M(x) N(x) = H(0).$$

Let  $Y = N(X)^{-1}X$ . Since  $N(X)$  is invertible, the map sending  $X$  to  $Y$  maps  $V$  diffeomorphically into its image. Also,

$$\begin{aligned} X^t M(X) X &= Y^t N(X)^t M(X) N(X) Y \\ &= Y^t H(0) Y \\ &= \sum_{2 \leq i \leq \lambda+1} Y_i^2 - \sum_{\lambda+2 \leq i \leq k} Y_i^2. \end{aligned}$$

□

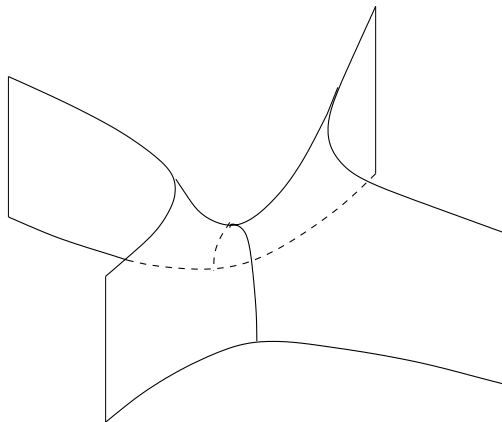
Using Lemma 7.15, we observe that in a small enough neighborhood of a critical point, a hypersurface behaves like one defined by a quadratic equation. So it suffices to analyze the change in the homotopy type of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq x}$  as  $x$  crosses 0 and the hypersurface defined by a quadratic polynomial of a very simple form. The change in the homotopy type consists in “attaching a handle along its boundary”, which is the process we describe now.

A  $j$ -ball is an embedding of  $\overline{B}_j(0, 1)$ , the closed  $j$ -dimensional ball with radius 1, in  $\text{Zer}(Q, \mathbb{R}^k)$ . It is a homeomorphic image of  $\overline{B}_j(0, 1)$  in  $\text{Zer}(Q, \mathbb{R}^k)$ .

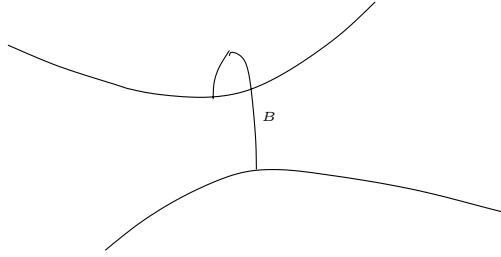
Let

$$P = X_1 - \sum_{2 \leq i \leq \lambda+1} X_i^2 + \sum_{\lambda+2 \leq i \leq k} X_i^2,$$

and  $\pi$  the projection onto the  $X_1$  axis restricted to  $\text{Zer}(P, \mathbb{R}^k)$ .



**Fig. 7.2.** The surface  $\text{Zer}(X_1 - X_2^2 + X_3^2, \mathbb{R}^3)$  near the origin



**Fig. 7.3.** The retract of  $\text{Zer}(X_1 - X_2^2 + X_3^2, \mathbb{R}^3)$  near the origin

Let  $B$  be the set defined by

$$X_2 = \dots = X_{\lambda+1} = 0, X_1 = -\sum_{\lambda+2 \leq i \leq k} X_i^2, -\epsilon \leq X_1 \leq 0.$$

Note that  $B$  is a  $(k - \lambda - 1)$ -ball and  $B \cap \text{Zer}(P, \mathbb{R}^k)_{\leq -\epsilon}$  is the set defined by

$$X_2 = \dots = X_{\lambda+1} = 0, X_1 = -\epsilon, \sum_{\lambda+2 \leq i \leq k} X_i^2 = \epsilon,$$

which is also the boundary of  $B$ .

**Lemma 7.16.** *For all sufficiently small  $\epsilon > 0$ , and  $r > 2\sqrt{\epsilon}$ , there exists a vector field  $\Gamma'$  on  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \setminus B$ , having the following properties:*

1. *Outside the ball  $B_k(r)$ ,  $2\epsilon\Gamma'$  equals the gradient vector field,  $\Gamma$ , of  $\pi$  on  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]}$ .*
2. *Associated to  $\Gamma'$  there is an one parameter continuous family of smooth maps  $\alpha_t: \text{Zer}(P, \mathbb{R}^k)_\epsilon \rightarrow \text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]}$ ,  $t \in [0, 1)$ , such that for  $x \in \text{Zer}(P, \mathbb{R}^k)_\epsilon$ ,  $t \in [0, 1)$ ,*
  - a) *Each  $\alpha_t$  is injective,*
  - b)  $\frac{d\alpha_t(x)}{dt} = \Gamma'(\alpha_t(x))$ ,
  - c)  $\alpha_0(x) = x$ ,
  - d)  $\lim_{t \rightarrow 1} \alpha_t(x) \in \text{Zer}(P, \mathbb{R}^k)_{-\epsilon} \cup B$ ,
  - e) *for every  $y \in \text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \setminus B$  there exists a unique  $z \in \text{Zer}(P, \mathbb{R}^k)_\epsilon$  and  $t \in [0, 1)$  such that  $\alpha_t(z) = y$ .*

**Proof of Lemma 7.16:** In the following, we consider  $\mathbb{R}^{k-1}$  as a product of the coordinate subspaces spanned by  $X_2, \dots, X_{\lambda+1}$  and  $X_{\lambda+2}, \dots, X_k$ , respectively, and denote by  $Y$  (resp.  $Z$ ) the vector of variables  $(X_2, \dots, X_{\lambda+1})$  (resp.  $(X_{\lambda+2}, \dots, X_k)$ ). We denote by  $\phi: \mathbb{R}^k \rightarrow \mathbb{R}^{k-1}$  the projection map onto the hyperplane  $X_1 = 0$ . Let  $S = \phi(\bar{B}_k(r))$ .

We depict the flow lines of the flow we are going to construct (projected onto the hyperplane defined by  $X_1 = 0$ ) in the case when  $k = 3$  and  $\lambda = 1$  in Figure 7.4.

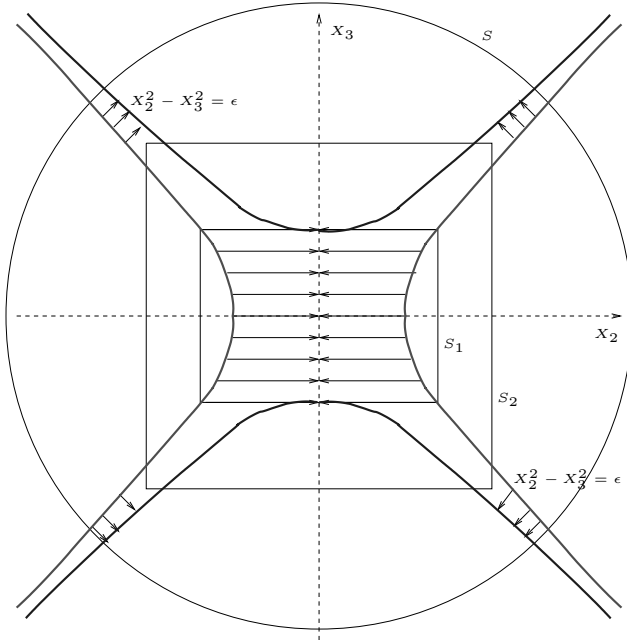


Fig. 7.4.  $S_1$  and  $S_2$

Consider the following two subsets of  $S$ .

$$S_1 = \overline{B}_\lambda(\sqrt{2}\epsilon) \times \overline{B}_{k-1-\lambda}(\sqrt{\epsilon})$$

and

$$S_2 = \overline{B}_\lambda(2\sqrt{\epsilon}) \times \overline{B}_{k-1-\lambda}(2\sqrt{\epsilon}).$$

In  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_1)$ , consider the flow lines whose projection onto the hyperplane  $X_1 = 0$  are straight segments joining the points  $(y_2, \dots, y_k) \in \phi(\text{Zer}(P, \mathbb{R}^k)_\epsilon)$  to  $(0, \dots, 0, y_{\lambda+2}, \dots, y_k)$ .

These correspond to the vector field on  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_1) \setminus B$  defined by

$$\Gamma_1 = \left( -\frac{1}{|Z|^2 + \epsilon}, \frac{-Y}{2|Y|^2(|Z|^2 + \epsilon)}, 0 \right).$$

Let  $p = (\epsilon, y, z) \in \text{Zer}(P, \mathbb{R}^k)_\epsilon \cap \phi^{-1}(S_1)$  and  $q$  the point in  $\text{Zer}(P, \mathbb{R}^k)$  having the same  $Z$  coordinates but having  $Y = 0$ . Then,  $\pi(q) = |z|^2 + \epsilon$ . Thus, the decreases uniformly from  $\epsilon$  to  $-|z|^2$  along the flow lines of the vector field  $\Gamma_1$ . For a point  $p = (x_1, y, z) \in \text{Zer}(Q, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_1) \setminus B$ , we denote by  $g(p)$  the limiting point on the flow line through  $p$  of the vector field  $\Gamma_1$  as it approaches  $Y = 0$ .

In  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S \setminus S_2)$ , consider the flow lines of the vector field

$$\Gamma_2 = \left( -\frac{1}{2\epsilon}, -\frac{Y}{4\epsilon(|Y|^2 + |Z|^2)}, \frac{Z}{4\epsilon(|Y|^2 + |Z|^2)} \right).$$

Notice that  $\Gamma_2$  is  $\frac{1}{2\epsilon}$  times the gradient vector field on

$$\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S \setminus S_2).$$

For a point  $p = (x_1, y, z) \in \text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S \setminus S_2)$ , we denote by  $g(p)$  the point on the flow line through  $p$  of the vector field  $\Gamma_2$  such that  $\pi(g(p)) = -\epsilon$ .

We patch these vector fields together in

$$\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_2 \setminus S_1)$$

using a  $C^\infty$  function that is 0 in  $S_1$  and 1 outside  $S_2$ . Such a function  $\mu: \mathbb{R}^{k-1} \rightarrow \mathbb{R}$  can be constructed as follows. Define

$$\lambda(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ 1 - 2^{-\frac{1}{4x^2}} & \text{if } 0 < x \leq \frac{1}{2}, \\ 2^{-\frac{1}{4(1-x)^2}} & \text{if } \frac{1}{2} < x \leq 1, \\ 1 & \text{if } x \geq 1. \end{cases}$$

Take

$$\mu(y, z) = \lambda\left(\frac{|y| - \sqrt{2}\epsilon}{\sqrt{2}\epsilon(\sqrt{2} - 1)}\right) \lambda\left(\frac{|z| - \sqrt{\epsilon}}{\sqrt{\epsilon}}\right).$$

Then, on  $\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_2 \setminus S_1)$  we consider the vector field

$$\Gamma'(p) = \mu(\phi(p))\Gamma_2(p) + (1 - \mu(\phi(p)))\Gamma_1(p).$$

Notice that it agrees with the vector fields defined on

$$\text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S \setminus S_2), \text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_1).$$

For a point  $p = (x_1, y, z) \in \text{Zer}(Q, \mathbb{R}^k)_{[-\epsilon, \epsilon]} \cap \phi^{-1}(S_2 \setminus S_1)$ , we denote by  $g(p)$  the point on the flow line through  $p$  of the vector field  $\Gamma_2$  such that  $\pi(g(p)) = -\epsilon$ .

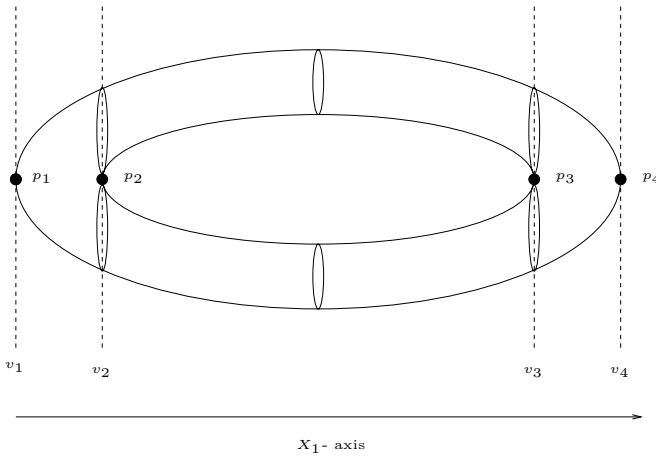
Denote the flow through a point  $p \in \text{Zer}(P, \mathbb{R}^k)_\epsilon \cap \phi^{-1}(S)$  of the vector field  $\Gamma'$  by  $\gamma_p: [0, 1] \rightarrow \text{Zer}(P, \mathbb{R}^k)_{[-\epsilon, \epsilon]}$ , with  $\gamma_p(0) = p$ .

For  $x \in \text{Zer}(P, \mathbb{R}^k)_\epsilon$  and  $t \in [0, 1]$ , define  $\alpha_t(x) = \gamma_x(t)$ . By construction of the vector field  $\Gamma$ ,  $\alpha_t$  has the required properties. □

Before proving Theorem 7.14 it is instructive to consider an example.



*Example 7.17.* Consider a smooth torus in  $\mathbb{R}^3$  (see Figure 7.5). There are four critical points  $p_1, p_2, p_3$  and  $p_4$  with critical values  $v_1, v_2, v_3$  and  $v_4$  and indices 2, 1, 1 and 0 respectively, for the projection map to the  $X_1$  coordinate.



**Fig. 7.5.** Changes in the homotopy type of the smooth torus in  $\mathbb{R}^3$  at the critical values

The changes in homotopy type at the corresponding critical values are described as follows: At the critical value  $v_1$  we add a 0-dimensional ball. At the critical values  $v_2$  and  $v_3$  we add 1-dimensional balls and finally at  $v_4$  we add a 2-dimensional ball.  $\square$

**Proof of Theorem 7.14:** We construct a vector field  $\Gamma'$  on  $\text{Zer}(Q, \mathbb{R}^k)_{[c-\epsilon, c+\epsilon]}$  that agrees with the gradient vector field  $\Gamma$  everywhere except in a small neighborhood of the critical point  $p$ . At the critical point  $p$ , we use Lemma 7.15 to reduce to the quadratic case and then use Lemma 7.16 to construct a vector field in a neighborhood of the critical point that agrees with  $\Gamma$  outside the neighborhood. We now use this vector field, as in the proof of Theorem 7.5, to obtain the required homotopy equivalence.  $\square$

We also need to analyze the topological changes that occur to sets bounded by non-singular algebraic hypersurfaces.

We are also going to prove the following versions of Theorem 7.5 (Morse Lemma A) and Theorem 7.14 (Morse Lemma B).

**Proposition 7.18.** *Let  $S$  be a bounded set defined by  $Q \geq 0$ , bounded by the non-singular algebraic hypersurface  $\text{Zer}(Q, \mathbb{R}^k)$ . Let  $[a, b]$  be an interval containing no critical value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ . Then  $S_{[a,b]}$  is homeomorphic to  $S_a \times [a, b]$  and  $S_{\leq a}$  is homotopy equivalent to  $S_{\leq b}$ .*

**Proposition 7.19.** *Let  $S$  be a bounded set defined by  $Q \geq 0$ , bounded by the non-singular algebraic hypersurface  $\text{Zer}(Q, \mathbb{R}^k)$ . Suppose that the projection  $\pi$  to the  $X_1$ -axis is a Morse function. Let  $p$  be the non-degenerate critical point of  $\pi$  on  $\partial W$  of index  $\lambda$  such that  $\pi(p) = c$ . For all sufficiently small  $\epsilon > 0$ , the set  $S_{\leq c+\epsilon}$  has*

- the homotopy type of  $S_{\leq c-\epsilon}$  if  $(\partial Q/\partial X_1)(p) < 0$ ,
- the homotopy type of the union of  $S_{\leq c-\epsilon}$  with a ball of dimension  $k-1-\lambda$  attached along its boundary, if  $(\partial Q/\partial X_1)(p) > 0$ .

*Example 7.20.* Consider the set in  $\mathbb{R}^3$  bounded by the smooth torus. Suppose that this set is defined by the single inequality  $Q \geq 0$ . In other words,  $Q$  is positive in the interior of the torus and negative outside. Referring back to Figure 7.5, we see that at the critical points  $p_2$  and  $p_4$ ,  $(\partial Q/\partial X_1)(p) < 0$  and hence according to Proposition 7.19 there is no change in the homotopy type at the two corresponding critical values  $v_2$  and  $v_4$ . However,  $(\partial Q/\partial X_1)(p) > 0$  at  $p_1$  and  $p_3$  and hence we add a 0-dimensional and an 1-dimensional balls at the two critical values  $v_1$  and  $v_3$  respectively.  $\square$

**Proof of Proposition 7.18:** Suppose that  $S$ , defined by  $Q \geq 0$ , is bounded by the non-singular algebraic hypersurface  $\text{Zer}(Q, \mathbb{R}^k)$ . We introduce a new variable,  $X_{k+1}$ , and consider the polynomial  $Q_+ = Q - X_{k+1}^2$  and the corresponding algebraic set  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$ . Let  $\phi: \mathbb{R}^{k+1} \rightarrow \mathbb{R}^k$  be the projection map to the first  $k$  coordinates.

Topologically,  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  consists of two copies of  $S$  glued along  $\text{Zer}(Q, \mathbb{R}^k)$ . Moreover, denoting by  $\pi'$  the projection from  $\mathbb{R}^{k+1}$  to  $\mathbb{R}$  forgetting the last  $k$  coordinates,  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  is non-singular and the critical points of  $\pi'$  on  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  are the critical points of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  (considering  $\text{Zer}(Q, \mathbb{R}^k)$  as a subset of the hyperplane defined by the equation  $X_{k+1} = 0$ ). We denote by  $\Gamma_+$  the gradient vector field on  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$ .

Since  $Q_+$  is a polynomial in  $X_1, \dots, X_k$  and  $X_{k+1}^2$ , the gradient vector field  $\Gamma_+$  on  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  is symmetric with respect to the reflection changing  $X_{k+1}$  to  $-X_{k+1}$ . Hence, we can project  $\Gamma_+$  and its associated flowlines down to the hyperplane defined by  $X_{k+1} = 0$  and get a vector field as well as its flowlines in  $S$ .

Now, the proof is exactly the same as the proof of Theorem 7.5 above, using the vector field  $\Gamma_+$  instead of  $\Gamma$ , and projecting the associated vector field down to  $\mathbb{R}^k$ , noting that the critical values of the projection map onto the first coordinate restricted to  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  are the same as those of  $\text{Zer}(Q, \mathbb{R}^k)$ .  $\square$

For the proof of Proposition 7.19, we first study the quadratic case.

Let  $\pi$  the projection onto the  $X_1$  axis and

$$P = X_1 - \sum_{2 \leq i \leq \lambda+1} X_i^2 + \sum_{\lambda+2 \leq i \leq k} X_i^2.$$

Let  $B_+$  be the set defined by

$$X_2 = \dots = X_{\lambda+1} = 0, X_1 = - \sum_{\lambda+2 \leq i \leq k} X_i^2, -\epsilon \leq X_1 \leq 0,$$

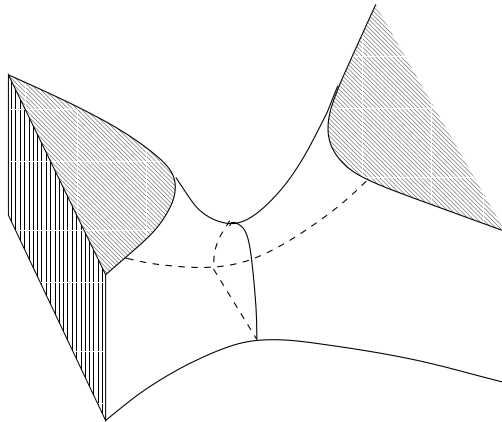
and let  $B_-$  be the set defined by

$$X_2 = \dots = X_{\lambda+1} = 0, X_1 \leq - \sum_{\lambda+2 \leq i \leq k} X_i^2, -\epsilon \leq X_1 \leq 0.$$

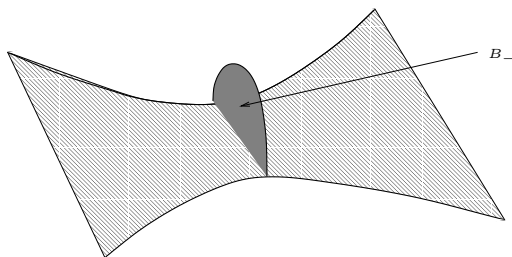
Note that,  $B_+$  is a  $(k - \lambda - 1)$ -ball and  $B_- \cap \text{Zer}(P, \mathbb{R}^k)_{\leq -\epsilon}$  is the set defined by

$$X_2 = \dots = X_{\lambda+1} = 0, X_1 = -\epsilon, \sum_{\lambda+2 \leq i \leq k} X_i^2 \leq \epsilon,$$

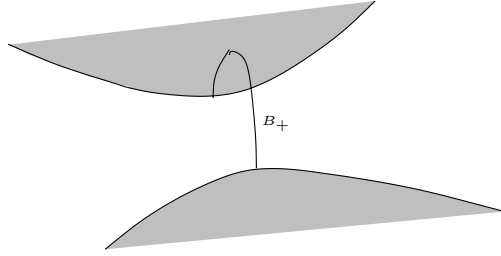
which is also the boundary of  $B_+$ .



**Fig. 7.6.** Set defined by  $X_1 - X_2^2 + X_3^2 \leq 0$  near the origin



**Fig. 7.7.** Retract of the set  $X_1 - X_2^2 + X_3^2 \leq 0$



**Fig. 7.8.** Retract of the set  $X_1 - X_2^2 + X_3^2 \geq 0$

**Lemma 7.21.** Let  $P_+ = P - X_{k+1}^2$ ,  $P_- = P + X_{k+1}^2$ .

1. Let  $S'$  be the set defined by  $P \geq 0$ . Then, for all sufficiently small  $\epsilon > 0$  and  $r > 2\sqrt{\epsilon}$ , there exists a vector field  $\Gamma'_+$  on  $S'_{[-\epsilon, \epsilon]} \setminus B_+$ , having the following properties:

a) Outside the ball  $B_k(r)$ ,  $2\epsilon\Gamma'_+$  equals the projection on  $\mathbb{R}^k$  of the gradient vector field,  $\Gamma_+$ , of  $\pi$  on  $\text{Zer}(P_+, \mathbb{R}^{k+1})_{[-\epsilon, \epsilon]}$ .

b) Associated to  $\Gamma'_+$ , there is a one parameter family of smooth maps  $\alpha_t^+ : S'_\epsilon \rightarrow S'_{[-\epsilon, \epsilon]}$ ,  $t \in [0, 1)$ , such that for  $x \in S'_\epsilon$ ,  $t \in [0, 1)$ ,

i. Each  $\alpha_t^+$  is injective,

ii.

$$\frac{d\alpha_t^+(x)}{dt} = \Gamma'_+(\alpha_t^+(x)),$$

iii.  $\alpha_0^+(x) = x$ ,

iv.  $\lim_{t \rightarrow 1} \alpha_t^+(x) \in S'_{-\epsilon} \cup B_+$  and,

v. for every  $y \in S'_{[-\epsilon, \epsilon]} \setminus B_+$  there exists a unique  $z \in S'_\epsilon$  and  $t \in [0, 1)$  such that  $\alpha_t(z) = y$ .

2. Similarly, let  $T'$  be the set defined by  $P \leq 0$ . Then, for all sufficiently small  $\epsilon > 0$  and  $r > 2\sqrt{\epsilon}$ , there exists a vector field  $\Gamma'_-$  on  $T'_{[-\epsilon, \epsilon]} \setminus B_+$  having the following properties:

a) Outside the ball  $B_k(r)$ ,  $2\epsilon\Gamma'_-$  the projection on  $\mathbb{R}^k$  of the gradient vector field,  $\Gamma_-$ , of  $\pi$  on  $\text{Zer}(P_-, \mathbb{R}^{k+1})_{[-\epsilon, \epsilon]}$ .

b) Associated to  $\Gamma'_-$ , there is a one parameter continuous family of smooth maps  $\alpha_t^- : T'_\epsilon \rightarrow T'_{[-\epsilon, \epsilon]}$ ,  $t \in [0, 1)$ , such that for  $x \in T'_\epsilon$ ,  $t \in [0, 1)$

i. Each  $\alpha_t^-$  is injective,

ii.

$$\frac{d\alpha_t^-(x)}{dt} = \Gamma'_-(\alpha_t^-(x)),$$

iii.  $\alpha_0^-(x) = x$ ,

iv.  $\lim_{t \rightarrow 1} \alpha_t^-(x) \in T'_{-\epsilon} \cup B_-$  and,

v. for every  $y \in T'_{[-\epsilon, \epsilon]} \setminus B_-$ , there exists a unique  $z \in T'_\epsilon$  and  $t \in [0, 1)$  such that  $\alpha_t(z) = y$ .

**Proof:** Since  $P_+$  (resp.  $P_-$ ) is a polynomial in  $X_1, \dots, X_k$  and  $X_{k+1}^2$ , the gradient vector field  $\Gamma_+$  (resp.  $\Gamma_-$ ) on  $\text{Zer}(P_+, \mathbb{R}^{k+1})$  (resp.  $\text{Zer}(P_-, \mathbb{R}^{k+1})$ ) is symmetric with respect to the reflection changing  $X_{k+1}$  to  $-X_{k+1}$ . Hence, we can project  $\Gamma_+$  (resp.  $\Gamma_-$ ) and its associated flowlines down to the hyperplane defined by  $X_{k+1} = 0$  and get a vector field  $\Gamma_+^*$  (resp.  $\Gamma_-^*$ ) as well as its flowlines in  $S'$  (resp.  $T'$ ).

1. Apply Lemma 7.16 to  $\text{Zer}(P_+, \mathbb{R}^k)$  to obtain a vector field  $\Gamma'_+$  on

$$\text{Zer}(P_+, \mathbb{R}^{k+1})_{[-\epsilon, \epsilon]} \setminus B_+$$

coinciding with  $\Gamma_+^*$ . Figure 7.8 illustrates the situation in the case  $k = 3$  and  $\lambda = 1$ .

2. Apply Lemma 7.16 to  $\text{Zer}(Q_-, \mathbb{R}^k)$  to obtain a vector field  $\Gamma'_-$  on

$$\text{Zer}(Q_-, \mathbb{R}^{k+1})_{[-\epsilon, \epsilon]} \setminus \phi^{-1}(B_-)$$

coinciding with  $\Gamma_-^*$ . Figures 7.8 and 7.8 illustrate the situation in the case  $k = 3$  and  $\lambda = 1$ . □

We are now in a position to prove Proposition 7.19.

**Proof of Proposition 7.19:** First, use Lemma 7.15 to reduce to the quadratic case, and then use Lemma 7.21, noting that the sign of  $\partial Q / \partial X_1 \} (p)$  determines which case we are in. □

## 7.2 Sum of the Betti Numbers of Real Algebraic Sets

For a closed semi-algebraic set  $S$ , let  $b(S)$  denote the sum of the Betti numbers of the simplicial homology groups of  $S$ . It follows from the definitions of Chapter 6 that  $b(S)$  is finite (see page 198).

According to Theorem 5.47, there are a finite number of algebraic subsets of  $\mathbb{R}^k$  defined by polynomials of degree at most  $d$ , say  $V_1, \dots, V_p$ , such that any algebraic subset  $V$  of  $\mathbb{R}^k$  so defined is semi-algebraically homeomorphic to one of the  $V_i$ . It follows immediately that any algebraic subset of  $\mathbb{R}^k$  defined by polynomials of degree at most  $d$  is such that  $b(V) \leq \max \{b(V_1), \dots, b(V_p)\}$ . Let  $b(k, d)$  be the smallest integer which bounds the sum of the Betti numbers of any algebraic set defined by polynomials of degree  $d$  in  $\mathbb{R}^k$ . The goal of this section is to bound the Betti numbers of a bounded non-singular algebraic hypersurface in terms of the number of critical values of a function defined on it and to obtain explicit bounds for  $b(k, d)$ .

*Remark 7.22.* Note that  $b(k, d) \geq d^k$  since the solutions to the system of equations,

$$(X_1 - 1)(X_1 - 2) \cdots (X_1 - d) = \cdots = (X_k - 1)(X_k - 2) \cdots (X_k - d) = 0$$

consist of  $d^k$  isolated points and the only non-zero Betti number of this set is  $b_0 = d^k$ . (Recall that for a closed and bounded semi-algebraic set  $S$ ,  $b_0(S)$  is the number of semi-algebraically connected components of  $S$  by Proposition 6.34.)  $\square$

We are going to prove the following theorem.

**Theorem 7.23. [Oleinik-Petrovski/Thom/Milnor bound]**

$$b(k, d) \leq d(2d - 1)^{k-1}.$$

The method for proving Theorem 7.23 will be to use Theorems 7.5 and 7.14, which give enough information about the homotopy type of  $\text{Zer}(Q, \mathbb{R}^k)$  to enable us to bound  $b(\text{Zer}(Q, \mathbb{R}^k))$  in terms of the number of critical points of  $\pi$ .

A first consequence of Theorems 7.5 and 7.14 is the following result.

**Theorem 7.24.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface such that the projection  $\pi$  on the  $X_1$ -axis is a Morse function. For  $0 \leq i \leq k-1$ , let  $c_i$  be the number of critical points of  $\pi$  restricted to  $\text{Zer}(Q, \mathbb{R}^k)$ , of index  $i$ . Then,*

$$\begin{aligned} b(\text{Zer}(Q, \mathbb{R}^k)) &\leq \sum_{\substack{i=0 \\ k-1}}^{k-1} c_i, \\ \chi(\text{Zer}(Q, \mathbb{R}^k)) &= \sum_{i=0}^{k-1} (-1)^{k-1-i} c_i. \end{aligned}$$

*In particular,  $b(\text{Zer}(Q, \mathbb{R}^k))$  is bounded by the number of critical points of  $\pi$  restricted to  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof:** Let  $v_1 < v_2 < \dots < v_\ell$  be the critical values of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  and  $p_i$  the corresponding critical points, such that  $\pi(p_i) = v_i$ . Let  $\lambda_i$  be the index of the critical point  $p_i$ . We first prove that  $b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i}) \leq i$ .

First note that  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_1}$  is  $\{p_1\}$  and hence

$$b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_1}) = b_0(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_1}) = 1.$$

By Theorem 7.5, the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i+1}-\epsilon}$  is homotopy equivalent to the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}$  for any small enough  $\epsilon > 0$ , and thus

$$b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i+1}-\epsilon}) = b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}).$$

By Theorem 7.14, the homotopy type of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}$  is that of the union of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i-\epsilon}$  with a topological ball. Recall from Proposition 6.44 that if  $S_1, S_2$  are two closed semi-algebraic sets with non-empty intersection, then

$$b_i(S_1 \cup S_2) \leq b_i(S_1) + b_i(S_2) + b_{i-1}(S_1 \cap S_2), \quad 0 \leq i \leq k-1.$$

Recall also from Proposition 6.34 that for a closed and bounded semi-algebraic set  $S$ ,  $b_0(S)$  equals the number of connected components of  $S$ . Since,  $S_1 \cap S_2 \neq \emptyset$ , for  $i = 0$  we have the stronger inequality,

$$b_0(S_1 \cup S_2) \leq b_0(S_1) + b_0(S_2) - 1.$$

By Proposition 6.37, for  $\lambda > 1$  we have that

$$\begin{aligned} b_0(B_\lambda) &= b_0(S^{\lambda-1}) \\ &= b_{\lambda-1}(S^{\lambda-1}) \\ &= 1, \\ b_i(B_\lambda) &= 0, i > 0, \\ b_i(S^{\lambda-1}) &= 0, 0 < i < \lambda - 1. \end{aligned}$$

It follows that, for  $\lambda > 1$ , attaching a  $\lambda$ -ball can increase  $b_\lambda$  by at most one, and none of the other Betti numbers can increase.

For  $\lambda = 1$ ,  $b_{\lambda-1}(S^{\lambda-1}) = b_0(S^0) = 2$ . It is an exercise to show that in this case,  $b_1$  can increase by at most one and no other Betti numbers can increase. (Hint. The number of cycles in a graph can increase by at most one on addition of an edge.)

It thus follows that

$$b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}) \leq b(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i-\epsilon}) + 1.$$

This proves the first part of the lemma.

We next prove that for  $1 < i \leq \ell$  and small enough  $\epsilon > 0$ ,

$$\chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}) = \chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i-1}+\epsilon}) + (-1)^{k-1-\lambda_i}.$$

By Theorem 7.5, the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i-\epsilon}$  is homotopy equivalent to the set  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i-1}+\epsilon}$  for any small enough  $\epsilon > 0$ , and thus

$$\chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i-\epsilon}) = \chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i-1}+\epsilon}).$$

By Theorem 7.14, the homotopy type of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}$  is that of the union of  $\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i-\epsilon}$  with a topological ball of dimension  $k - 1 - \lambda_i$ . Recall from Corollary 6.36 (Equation 6.36) that if  $S_1, S_2$  are two closed and bounded semi-algebraic sets with non-empty intersection, then

$$\chi(S_1 \cup S_2) = \chi(S_1) + \chi(S_2) - \chi(S_1 \cap S_2).$$

Hence,

$$\begin{aligned} \chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i+\epsilon}) &= \chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i-1}+\epsilon}) \\ &= \chi(\overline{B}_{k-1-\lambda_i}) \\ &\quad - \chi(S^{k-2-\lambda_i}). \end{aligned}$$

Now, it follows from Proposition 6.37 and the definition of Euler-Poincaré characteristic, that  $\chi(\overline{B}_{k-1-\lambda_i}) = 1$  and  $\chi(S^{k-2-\lambda_i}) = 1 + (-1)^{k-2-\lambda_i}$ .

Substituting in the equation above we obtain that

$$\chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_i + \epsilon}) = \chi(\text{Zer}(Q, \mathbb{R}^k)_{\leq v_{i-1} + \epsilon}) + (-1)^{k-1-\lambda_i}.$$

The second part of the theorem is now an easy consequence.  $\square$

We shall need the slightly more general result.

**Proposition 7.25.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface such that the projection  $\pi$  on the  $X_1$ -axis has non-degenerate critical points on  $\text{Zer}(Q, \mathbb{R}^k)$ . For  $0 \leq i \leq k-1$ , let  $c_i$  be the number of critical points of  $\pi$  restricted to  $\text{Zer}(Q, \mathbb{R}^k)$ , of index  $i$ . Then,*

$$\begin{aligned} b(\text{Zer}(Q, \mathbb{R}^k)) &\leq \sum_{i=0}^{k-1} c_i, \\ \chi(\text{Zer}(Q, \mathbb{R}^k)) &= \sum_{i=0}^{k-1} (-1)^{k-1-i} c_i. \end{aligned}$$

In particular,  $b(\text{Zer}(Q, \mathbb{R}^k))$  is bounded by the number of critical points of  $\pi$  restricted to  $\text{Zer}(Q, \mathbb{R}^k)$ .

**Proof:** Use Lemma 7.13 and Theorem 7.24.  $\square$

Using Theorem 7.24, we can estimate the sum of the Betti numbers in the bounded case.

**Proposition 7.26.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface with  $Q$  a polynomial of degree  $d$ . Then*

$$b(\text{Zer}(Q, \mathbb{R}^k)) \leq d(d-1)^{k-1}.$$

**Proof:** Using Proposition 7.9, we can suppose that  $\pi$  is a Morse function. Applying Theorem 7.24 to the function  $\pi: \text{Zer}(Q, \mathbb{R}^k) \rightarrow \mathbb{R}$ , it follows that the sum of the Betti numbers of  $\text{Zer}(Q, \mathbb{R}^k)$  is less than or equal to the number of critical points of  $\pi$ . Now apply Proposition 7.10.  $\square$

In order to obtain Theorem 7.23, we will need the following Proposition.

**Proposition 7.27.** *Let  $S$  be a bounded set defined by  $Q \geq 0$ , bounded by the non-singular algebraic hypersurface  $\text{Zer}(Q, \mathbb{R}^k)$ . Let the projection map  $\pi$  be a Morse function on  $\text{Zer}(Q, \mathbb{R}^k)$ . Then, the sum of the Betti numbers of  $S$  is bounded by half the number of critical points of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof:** We use the notation of the proof of Proposition 7.18. Let  $v_1 < v_2 < \dots < v_\ell$  be the critical values of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  and  $p_1, \dots, p_\ell$  the corresponding critical points, such that  $\pi(p_i) = v_i$ . We denote by  $J$  the subset of  $\{1, \dots, \ell\}$  such that the direction of  $\text{Grad}(Q)(p)$  belongs to  $S$  (see Proposition 7.18).



We are going to prove that

$$b(S_{\leq v_i}) \leq \#(j \in J, j \leq i).$$

First note that  $S_{\leq v_1}$  is  $\{p_1\}$  and hence  $b(S_{\leq v_1}) = 1$ . By Proposition 7.18  $S_{\leq v_{i+1}-\epsilon}$  is homotopic to  $S_{\leq v_i+\epsilon}$  for any small enough  $\epsilon > 0$ , and thus

$$b(S_{\leq v_{i+1}-\epsilon}) = b(S_{\leq v_i+\epsilon}).$$

By Theorem 7.14, the homotopy type of  $S_{\leq v_i+\epsilon}$  is that of  $S_{\leq v_i-\epsilon}$  if  $i \notin J$  and that of the union of  $S_{\leq v_i-\epsilon}$  with a topological ball if  $i \in J$ .

It follows that

$$\begin{cases} b(S_{\leq v_i+\epsilon}) = b(S_{\leq v_i-\epsilon}) & \text{if } i \notin J \\ b(S_{\leq v_i+\epsilon}) \leq b(S_{\leq v_i-\epsilon}) + 1 & \text{if } i \in J. \end{cases}$$

By switching the direction of the  $X_1$  axis if necessary, we can always ensure that  $\#(J)$  is at most half of the critical points.  $\square$

**Proposition 7.28.** *If  $\mathbb{R} = \mathbb{R}$ ,*

$$b(k, d) \leq d(2d - 1)^{k-1}.$$

**Proof:** Let  $V = \text{Zer}(\{P_1, \dots, P_\ell\}, \mathbb{R}^k)$  with the the degrees of the  $P_i$ 's bounded by  $d$ . By remark on page 226, it suffices to estimate the sum of the Betti numbers of  $V \cap \bar{B}_k(0, r)$ . Let

$$F(X) = \frac{P_1^2 + \dots + P_\ell^2}{r^2 - \|X\|^2}.$$

By Sard's theorem (Theorem 5.56), the set of critical values of  $F$  is finite. Hence, there is a positive  $a \in \mathbb{R}$  so that no  $b \in (0, a)$  is a critical value of  $F$  and thus the set  $W_b = \{x \in \mathbb{R}^k \mid P(x, b) = 0\}$ , where

$$P(X, b) = P_1^2 + \dots + P_\ell^2 + b(\|X\|^2 - r^2)$$

is a non-singular hypersurface in  $\mathbb{R}^k$ . To see this observe that, for  $x \in \mathbb{R}^k$

$$P(x, b) = \partial P / \partial X_1(x, b) = \dots = \partial P / \partial X_k(x, b) = 0$$

implies that  $F(x) = b$  and  $\partial F / \partial X_1(x) = \dots = \partial F / \partial X_k(x) = 0$  implying that  $b$  is a critical value of  $F$  which is a contradiction.

Moreover,  $W_b$  is the boundary of the closed and bounded set

$$K_b = \{x \in \mathbb{R}^k \mid P(x, b) \leq 0\}.$$

By Proposition 7.26, the sum of the Betti numbers of  $W_b$  is less than or equal to  $2d(2d - 1)^{k-1}$ .

Also, using Proposition 7.27, the sum of the Betti numbers of  $K_b$  is at most half that of  $W_b$ .

We now claim that  $V \cap \overline{B}_k(0, r)$  is homotopy equivalent to  $K_b$  for all small enough  $b > 0$ . We replace  $b$  in the definition of the set  $K_b$  by a new variable  $T$ , and consider the set  $K \subset \mathbb{R}^{k+1}$  defined by  $\{(x, t) \in \mathbb{R}^{k+1} \mid P(x, t) \leq 0\}$ . Let  $\pi_X$  (resp.  $\pi_T$ ) denote the projection map onto the  $X$  (resp.  $T$ ) coordinates.

Clearly,  $V \cap \overline{B}_k(0, r) \subset K_b$ . By Theorem 5.46 (Semi-algebraic triviality), for all small enough  $b > 0$ , there exists a semi-algebraic homeomorphism,

$$\phi: K_b \times (0, b] \rightarrow K \cap \pi_T^{-1}((0, b]),$$

such that  $\pi_T(\phi(x, s)) = s$  and  $\phi$  is a semi-algebraic homeomorphism from  $V \cap B_k(0, r) \times (0, b]$  to itself.

Let  $G: K_b \times [0, b] \rightarrow K_b$  be the map defined by  $G(x, s) = \pi_X(\phi(x, s))$  for  $s > 0$  and  $G(x, 0) = \lim_{s \rightarrow 0^+} \pi_X(\phi(x, s))$ . Let  $g: K_b \rightarrow V \cap \overline{B}_k(0, r)$  be the map  $G(x, 0)$  and  $i: V \cap \overline{B}_k(0, r) \rightarrow K_b$  the inclusion map. Using the homotopy  $G$ , we see that  $i \circ g \sim \text{Id}_{K_b}$ , and  $g \circ i \sim \text{Id}_{V \cap B_k(0, r)}$ , which shows that  $V \cap B_k(0, r)$  is homotopy equivalent to  $K_b$  as claimed.

Hence,

$$b(V \cap \overline{B}_k(0, r)) = b(K_b) \leq 1/2 b(W_b) \leq d(2d-1)^{k-1}. \quad \square$$

**Proof of Theorem 7.23:** It only remains to prove that Proposition 7.28 is valid for any real closed field  $\mathbb{R}$ . We first work over the field of real algebraic numbers  $\mathbb{R}_{\text{alg}}$ . We identify a system of  $\ell$  polynomials  $(P_1, \dots, P_\ell)$  in  $k$  variables of degree less than or equal to  $d$  with the point of  $\mathbb{R}_{\text{alg}}^N$ ,  $N = \ell \binom{k+d-1}{d}$ , whose coordinates are the coefficients of  $P_1, \dots, P_\ell$ . Let

$$Z = \{(P_1, \dots, P_\ell, x) \in \mathbb{R}_{\text{alg}}^N \times \mathbb{R}_{\text{alg}}^k \mid P_1(x) = \dots = P_\ell(x) = 0\},$$

and let  $\Pi: Z \rightarrow \mathbb{R}_{\text{alg}}^N$  be the canonical projection. By Theorem 5.46 (Semi-algebraic Triviality), there exists a finite partition of  $\mathbb{R}_{\text{alg}}^N$  into semi-algebraic sets  $A_1, \dots, A_m$ , semi-algebraic sets  $F_1, \dots, F_m$  contained in  $\mathbb{R}_{\text{alg}}^k$ , and semi-algebraic homeomorphisms  $\theta_i: \Pi^{-1}(A_i) \rightarrow A_i \times F_i$ , for  $i = 1, \dots, m$ , such that the composition of  $\theta_i$  with the projection  $A_i \times F_i \rightarrow A_i$  is  $\Pi|_{\Pi^{-1}(A_i)}$ . The  $F_i$  are algebraic subsets of  $\mathbb{R}_{\text{alg}}^k$  defined by  $\ell$  equations of degree less than or equal to  $d$ . The sum of the Betti numbers of  $\text{Ext}(F_i, \mathbb{R})$  is less than or equal to  $d(2d-1)^{k-1}$ . So, by invariance of the homology groups under extension of real closed field (Section 6.2), the same bound holds for the sum of the Betti numbers of  $F_i$ . Now, let  $V \subset \mathbb{R}^k$  be defined by  $k$  equations  $P_1 = \dots = P_\ell = 0$  of degree less than or equal to  $d$  with coefficients in  $\mathbb{R}$ . We have

$$\text{Ext}(\Pi^{-1}, \mathbb{R})(P_1, \dots, P_\ell) = \{(P_1, \dots, P_\ell)\} \times V.$$

The point  $(P_1, \dots, P_\ell) \in \mathbb{R}^N$  belongs to some  $\text{Ext}(A_i, \mathbb{R})$ , and the semi-algebraic homeomorphism  $\text{Ext}(\theta_i, \mathbb{R})$  induces a semi-algebraic homeomorphism from  $V$  onto  $\text{Ext}(F_i, \mathbb{R})$ . Again, the sum of the Betti numbers of  $\text{Ext}(F_i, \mathbb{R})$  is less than or equal to  $d(2d-1)^{k-1}$ , and the same bound holds for the sum of the Betti numbers of  $V$ .  $\square$

### 7.3 Bounding the Betti Numbers of Realizations of Sign Conditions

Throughout this section, let  $\mathcal{Q}$  and  $\mathcal{P} \neq \emptyset$  be finite subsets of  $\mathbb{R}[X_1, \dots, X_k]$ , let  $Z = \text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ , and let  $k'$  be the dimension of  $Z = \text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ .

**Notation 7.29. [Realizable sign conditions]** We denote by

$$\text{SIGN}(\mathcal{P}) \subset \{0, 1, -1\}^{\mathcal{P}}$$

the set of all realizable sign conditions for  $\mathcal{P}$  over  $\mathbb{R}^k$ , and by

$$\text{SIGN}(\mathcal{P}, \mathcal{Q}) \subset \{0, 1, -1\}^{\mathcal{P}}$$

the set of all realizable sign conditions for  $\mathcal{P}$  over  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ . □

For  $\sigma \in \text{SIGN}(\mathcal{P}, \mathcal{Q})$ , let  $b_i(\sigma)$  denote the  $i$ -th Betti number of

$$\text{Reali}(\sigma, Z) = \{x \in \mathbb{R}^k \mid \bigwedge_{Q \in \mathcal{Q}} Q(x) = 0, \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma(P)\}.$$

Let  $b_i(\mathcal{Q}, \mathcal{P}) = \sum_{\sigma} b_i(\sigma)$ . Note that  $b_0(\mathcal{Q}, \mathcal{P})$  is the number of semi-algebraically connected components of basic semi-algebraic sets defined by  $\mathcal{P}$  over  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ .

We denote by  $\text{deg}(\mathcal{Q})$  the maximum of the degrees of the polynomials in  $\mathcal{Q}$  and write  $b_i(d, k, k', s)$  for the maximum of  $b_i(\mathcal{Q}, \mathcal{P})$  over all  $\mathcal{Q}, \mathcal{P}$ , where  $\mathcal{Q}$  and  $\mathcal{P}$  are finite subsets of  $\mathbb{R}[X_1, \dots, X_k]$ ,  $\text{deg}(\mathcal{Q}, \mathcal{P}) \leq d$  whose elements have degree at most  $d$ ,  $\#(\mathcal{P}) = s$  (i.e.  $\mathcal{P}$  has  $s$  elements), and the algebraic set  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$  has dimension  $k'$ .

**Theorem 7.30.**

$$b_i(d, k, k', s) \leq \sum_{1 \leq j \leq k' - i} \binom{s}{j} 4^j d (2d - 1)^{k-1}.$$

So we get, in particular a bound on the total number of semi-algebraically connected components of realizable sign conditions.

**Proposition 7.31.**

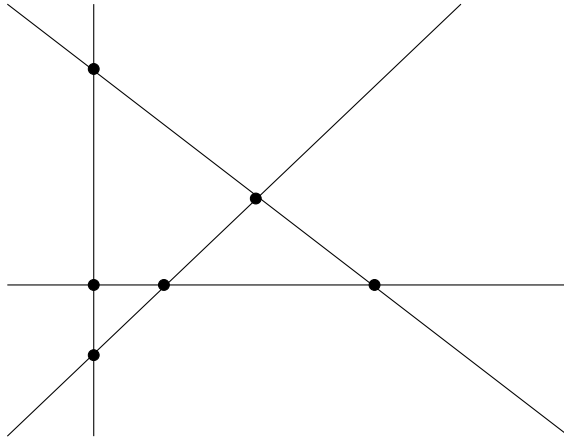
$$b_0(d, k, k', s) \leq \sum_{1 \leq j \leq k'} \binom{s}{j} 4^j d (2d - 1)^{k-1}.$$

*Remark 7.32.* When  $d = 1$ , i.e. when all equations are linear, it is easy to find directly a bound on the number of non-empty sign conditions. The number of non-empty sign conditions  $f(k', s)$  defined by  $s$  linear equations on a flat of dimension  $k'$  satisfies the recurrence relation

$$f(k', s + 1) \leq f(k', s) + 2 f(k' - 1, s),$$

since a flat  $L$  of dimension  $k' - 1$  meets at most  $f(k' - 1, s)$  non-empty sign condition defined by  $s$  polynomials on a flat of dimension  $k'$ , and each such non-empty sign condition is divided in at most three pieces by  $L$ .

In Figure 7.9 we depict the situation with four lines in  $\mathbb{R}^2$  defined by four linear polynomials. The number of realizable sign conditions in this case is easily seen to be 33.



**Fig. 7.9.** Four lines in  $\mathbb{R}^2$

Moreover, when the linear equations are in general position,

$$f(k', s + 1) = f(k', s) + 2 f(k' - 1, s). \tag{7.4}$$

Since  $f(k', 0) = 1$ , the solution to Equation (7.4) is given by

$$f(k', s) = \sum_{i=0}^{k'} \sum_{j=0}^{k'-i} \binom{s}{i} \binom{s-i}{j}. \tag{7.5}$$

Since all the realizations are convex and hence contractible, this bound on the number of non-empty sign conditions is also a bound on

$$b_0(1, k, k', s) = b(1, k', k', s)$$

We note that

$$f(k', s) \leq \sum_{1 \leq j \leq k'} \binom{s}{j} 4^j,$$

the right hand side being the bound appearing in Proposition 7.31 with  $d=1$ . □

The following proposition, Proposition 7.33, plays a key role in the proofs of these theorems. Part (a) of the proposition bounds the Betti numbers of a union of  $s$  semi-algebraic sets in  $\mathbb{R}^k$  in terms of the Betti numbers of the intersections of the sets taken at most  $k$  at a time.

Part (b) of the proposition is a dual version of Part (a) with unions being replaced by intersections and vice-versa, with an additional complication arising from the fact that the empty intersection, corresponding to the base case of the induction, is an arbitrary real algebraic variety of dimension  $k'$ , and is generally not acyclic.

Let  $S_1, \dots, S_s \subset \mathbb{R}^k$ ,  $s \geq 1$ , be closed semi-algebraic sets contained in a closed semi-algebraic set  $T$  of dimension  $k'$ . For  $1 \leq t \leq s$ , let  $S_{\leq t} = \bigcap_{1 \leq j \leq t} S_j$ , and  $S^{\leq t} = \bigcup_{1 \leq j \leq t} S_j$ . Also, for  $J \subset \{1, \dots, s\}$ ,  $J \neq \emptyset$ , let  $S_J = \bigcap_{j \in J} S_j$ , and  $S^J = \bigcup_{j \in J} S_j$ . Finally, let  $S^\emptyset = T$ .

**Proposition 7.33.**

a) For  $0 \leq i \leq k'$ ,

$$b_i(S^{\leq s}) \leq \sum_{j=1}^{i+1} \sum_{\substack{J \subset \{1, \dots, s\} \\ \#(J)=j}} b_{i-j+1}(S_J). \tag{7.6}$$

b) For  $0 \leq i \leq k'$ ,

$$b_i(S_{\leq s}) \leq \sum_{j=1}^{k'-i} \sum_{\substack{J \subset \{1, \dots, s\} \\ \#(J)=j}} b_{i+j-1}(S^J) + \binom{s}{k'-i} b_{k'}(S^\emptyset). \tag{7.7}$$

**Proof :** a) We prove the claim by induction on  $s$ . The statement is clearly true for  $s=1$ , since  $b_i(S_1)$  appears on the right hand side for  $j=1$  and  $J=\{1\}$ .

Using Proposition 6.44 (6.44), we have that

$$b_i(S^{\leq s}) \leq b_i(S^{\leq s-1}) + b_i(S_s) + b_{i-1}(S^{\leq s-1} \cap S_s). \tag{7.8}$$

Applying the induction hypothesis to the set  $S^{\leq s-1}$ , we deduce that

$$b_i(S^{\leq s-1}) \leq \sum_{j=1}^{i+1} \sum_{\substack{J \subset \{1, \dots, s-1\} \\ \#(J)=j}} b_{i-j+1}(S_J). \tag{7.9}$$

Next, we apply the induction hypothesis to the set

$$S^{\leq s-1} \cap S_s = \bigcup_{1 \leq j \leq s-1} (S_j \cap S_s)$$

to get that

$$b_{i-1}(S^{\leq s-1} \cap S_s) \leq \sum_{j=1}^i \sum_{\substack{J \subset \{1, \dots, s-1\} \\ \#(J)=j}} b_{i-j}(S_{J \cup \{s\}}). \tag{7.10}$$

Adding the inequalities (7.9) and (7.10), we get

$$b_i(S^{\leq s-1}) + b_i(S_s) + b_{i-1}(S^{\leq s-1} \cap S_s) \leq \sum_{j=1}^{i+1} \sum_{\substack{J \subset \{1, \dots, s\} \\ \#(J)=j}} b_{i-j+1}(S_J).$$

We conclude using (7.8).

b) We first prove the claim when  $s = 1$ . If  $0 \leq i \leq k' - 1$ , the claim is

$$b_i(S_1) \leq b_{k'}(S^\emptyset) + (b_i(S_1) + b_{k'}(S^\emptyset)),$$

which is clear. If  $i = k'$ , the claim is  $b_{k'}(S_1) \leq b_{k'}(S^\emptyset)$ . If the dimension of  $S_1$  is  $k'$ , consider the closure  $V$  of the complement of  $S_1$  in  $T$ . The intersection  $W$  of  $V$  with  $S_1$ , which is the boundary of  $S_1$ , has dimension strictly smaller than  $k'$  by Theorem 5.42 thus  $b_{k'}(W) = 0$ . Using Proposition 6.44

$$b_{k'}(S_1) + b_{k'}(V) \leq b_{k'}(S^\emptyset) + b_{k'}(W),$$

and the claim follows. On the other hand, if the dimension of  $S_1$  is strictly smaller than  $k'$ ,  $b_{k'}(S_1) = 0$ .

The claim is now proved by induction on  $s$ . Assume that the induction hypothesis (7.7) holds for  $s - 1$  and for all  $0 \leq i \leq k'$ .

From Proposition 6.44 (6.44), we have

$$b_i(S_{\leq s}) \leq b_i(S_{\leq s-1}) + b_i(S_s) + b_{i+1}(S_{\leq s-1} \cup S_s). \tag{7.11}$$

Applying the induction hypothesis to the set  $S_{\leq s-1}$ , we deduce that

$$\begin{aligned} b_i(S_{\leq s-1}) &\leq \sum_{j=1}^{k'-i} \sum_{\substack{J \subset \{1, \dots, s-1\} \\ \#(J)=j}} b_{i+j-1}(S^J) \\ &\quad + \binom{s-1}{k'-i} b_{k'}(S^\emptyset). \end{aligned}$$

Next, applying the induction hypothesis to the set,

$$S_{\leq s-1} \cup S_s = \bigcap_{1 \leq j \leq s-1} (S_j \cup S_s),$$

we get that

$$\begin{aligned} b_{i+1}(S_{\leq s-1} \cup S_s) &\leq \sum_{j=1}^{k'-i-1} \sum_{\substack{J \subset \{1, \dots, s-1\} \\ \#(J)=j}} b_{i+j}(S^{J \cup \{s\}}) \\ &\quad + \binom{s-1}{k'-i-1} b_{k'}(S^\emptyset). \end{aligned} \tag{7.12}$$

Adding the inequalities (7.11) and (7.12), we get

$$b_i(S_{\leq s}) \leq \sum_{j=1}^{k'-i} \sum_{\substack{J \subset \{1, \dots, s\} \\ \#(J)=j}} b_{i+j-1}(S^J) + \binom{s}{k'-i} b_{k'}(S^\emptyset).$$

We conclude using (7.11). □

Let  $\mathcal{P} = \{P_1, \dots, P_s\}$ , and let  $\delta$  be a new variable. We will consider the field  $\mathbb{R}(\delta)$  of algebraic Puiseux series in  $\delta$ , in which  $\delta$  is an infinitesimal.

Let  $S_i = \text{Reali}(P_i^2(P_i^2 - \delta^2) \geq 0, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle))$ ,  $1 \leq i \leq s$ , and let  $S$  be the intersection of the  $S_i$  with the closed ball in  $\mathbb{R}\langle\delta\rangle^k$  defined by

$$\delta^2 \left( \sum_{1 \leq i \leq k} X_i^2 \right) \leq 1.$$

In order to estimate  $b_i(S)$ , we prove that  $b_i(\mathcal{P}, \mathcal{Q})$  and  $b_i(S)$  are equal and we estimate  $b_i(S)$ .

**Proposition 7.34.**

$$b_i(\mathcal{P}, \mathcal{Q}) = b_i(S).$$

**Proof:** Consider a sign condition  $\sigma$  on  $\mathcal{P}$  such that, without loss of generality,

$$\begin{aligned} \sigma(P_i) &= 0 && \text{if } i \in I \\ \sigma(P_j) &= 1 && \text{if } j \in J \\ \sigma(P_\ell) &= -1 && \text{if } \ell \in \{1, \dots, s\} \setminus (I \cup J), \end{aligned}$$

and denote by  $\overline{\text{Reali}}(\sigma)$  the subset of  $\text{Ext}(Z, \mathbb{R}\langle\delta\rangle)$  defined by

$$\delta^2 \left( \sum_{1 \leq i \leq k} X_i^2 \right) \leq 1, P_i = 0, i \in I,$$

$$P_j \geq \delta, j \in J, P_\ell \leq -\delta, \ell \in \{1, \dots, s\} \setminus (I \cup J).$$

Note that  $S$  is the disjoint union of the  $\overline{\text{Reali}}(\sigma)$  for all realizable sign conditions  $\sigma$ .

Moreover, by definition of the homology groups of sign conditions (Notation 6.46)  $b_i(\sigma) = b_i(\overline{\text{Reali}}(\sigma))$ , so that

$$b_i(\mathcal{P}, \mathcal{Q}) = \sum_{\sigma} b_i(\sigma) = b_i(S). \quad \square$$

**Proposition 7.35.**

$$b_i(S) \leq \sum_{j=1}^{k'-i} \binom{s}{j} 4^j d (2d - 1)^{k-1}.$$

Before estimating  $b_i(S)$ , we estimate the Betti numbers of the following sets.

Let  $j \geq 1$ ,

$$V_j = \text{Reali} \left( \bigvee_{1 \leq i \leq j} P_i^2(P_i^2 - \delta^2) = 0, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle) \right),$$

and

$$W_j = \text{Reali} \left( \bigvee_{1 \leq i \leq j} P_i^2(P_i^2 - \delta^2) \geq 0, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle) \right).$$

Note that  $W_j$  is the union of  $S_1, \dots, S_j$ .

**Lemma 7.36.**

$$b_i(V_j) \leq (4^j - 1) d (2d - 1)^{k-1}.$$

**Proof:** Each of the sets

$$\text{Reali}(P_i^2(P_i^2 - \delta^2)) = 0, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle)$$

is the disjoint union of three algebraic sets, namely

$$\text{Reali}(P_i = 0, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle)),$$

$$\text{Reali}(P_i = \delta, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle)),$$

$$\text{Reali}(P_i = -\delta, \text{Ext}(Z, \mathbb{R}\langle\delta\rangle)).$$

Moreover, each Betti number of their union is bounded by the sum of the Betti numbers of all possible non-empty sets that can be obtained by taking, for  $1 \leq \ell \leq j$ ,  $\ell$ -ary intersections of these algebraic sets, using part (a) of Proposition 7.33. The number of possible  $\ell$ -ary intersection is  $\binom{j}{\ell}$ . Each such intersection is a disjoint union of  $3^\ell$  algebraic sets. The sum of the Betti numbers of each of these algebraic sets is bounded by  $d(2d - 1)^{k-1}$  by using Theorem 7.23.

Thus,

$$b_i(V_j) \leq \sum_{\ell=1}^j \binom{j}{\ell} 3^\ell d (2d - 1)^{k-1} = (4^j - 1) d (2d - 1)^{k-1}. \quad \square$$

**Lemma 7.37.**

$$b_i(W_j) \leq (4^j - 1) d (2d - 1)^{k-1} + b_i(Z).$$

**Proof:** Let  $Q_i = P_i^2(P_i^2 - \delta^2)$  and

$$F = \text{Reali}\left(\bigwedge_{1 \leq i \leq j} (Q_i \leq 0) \vee \bigvee_{1 \leq i \leq j} (Q_i = 0), \text{Ext}(Z, \mathbb{R}\langle\delta\rangle)\right).$$

Apply inequality (6.44), noting that

$$W_j \cup F = \text{Ext}(Z, \mathbb{R}\langle\delta\rangle), \quad W_j \cap F = W_{j,0}.$$

Since  $b_i(Z) = b_i(\text{Ext}(Z, \mathbb{R}\langle\delta\rangle))$ , we get that

$$b_i(W_j) \leq b_i(W_j \cap F) + b_i(W_j \cup F) = b_i(V_j) + b_i(Z).$$

We conclude using Lemma 7.36. □



**Proof of Proposition 7.35:** Using part b) of Proposition 7.33 and Lemma 7.37, we get

$$b_i(S) \leq \sum_{j=1}^{k'-i} \binom{s}{j} ((4^j - 1) d (2d - 1)^{k-1} + b_i(Z)) + \binom{s}{k'-i} b_{k'}(Z).$$

By Theorem 7.23, for all  $i < k'$ ,

$$b_i(Z) + b_{k'}(Z) \leq d(2d - 1)^{k-1}.$$

Thus, we have

$$b_i(S) \leq \sum_{j=1}^{k'-i} \binom{s}{j} 4^j d (2d - 1)^{k-1}. \quad \square$$

Theorem 7.30 follows clearly from Proposition 7.34 and Proposition 7.35.

### 7.4 Sum of the Betti Numbers of Closed Semi-algebraic Sets

Let  $\mathcal{P}$  and  $\mathcal{Q}$  be finite subsets of  $\mathbb{R}[X_1, \dots, X_k]$ .

A  $(\mathcal{Q}, \mathcal{P})$ -closed formula is a formula constructed as follows:

- For each  $P \in \mathcal{P}$ ,

$$\bigwedge_{Q \in \mathcal{Q}} Q = 0 \wedge P = 0, \quad \bigwedge_{Q \in \mathcal{Q}} Q = 0 \wedge P \geq 0, \quad \bigwedge_{Q \in \mathcal{Q}} Q = 0 \wedge P \leq 0,$$

- If  $\Phi_1$  and  $\Phi_2$  are  $(\mathcal{Q}, \mathcal{P})$ -closed formulas,  $\Phi_1 \wedge \Phi_2$  and  $\Phi_1 \vee \Phi_2$  are  $(\mathcal{Q}, \mathcal{P})$ -closed formulas.

Clearly,  $\text{Real}(\Phi)$ , the realization of a  $(\mathcal{Q}, \mathcal{P})$ -closed formula  $\Phi$ , is a closed semi-algebraic set. We denote by  $b(\Phi)$  the sum of its Betti numbers.

We write  $\bar{b}(d, k, k', s)$  for the maximum of  $b(\Phi)$ , where  $\Phi$  is a  $(\mathcal{Q}, \mathcal{P})$ -closed formula,  $\mathcal{Q}$  and  $\mathcal{P}$  are finite subsets of  $\mathbb{R}[X_1, \dots, X_k]$   $\text{deg}(\mathcal{Q}, \mathcal{P}) \leq d$ ,  $\#(\mathcal{P}) = s$ , and the algebraic set  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$  has dimension  $k'$ .

Our aim in this section is to prove the following result.

**Theorem 7.38.**

$$\bar{b}(d, k, k', s) \leq \sum_{i=0}^{k'} \sum_{j=1}^{k'-i} \binom{s}{j} 6^j d (2d - 1)^{k-1}.$$

For the proof of Theorem 7.38, we are going to introduce several infinitesimal quantities. Given a list of polynomials  $\mathcal{P} = \{P_1, \dots, P_s\}$  with coefficients in  $R$ , we introduce  $s$  new variables  $\delta_1, \dots, \delta_s$  and inductively define

$$\mathbb{R}\langle \delta_1, \dots, \delta_{i+1} \rangle = \mathbb{R}\langle \delta_1, \dots, \delta_i \rangle \langle \delta_{i+1} \rangle.$$

Note that  $\delta_{i+1}$  is infinitesimal with respect to  $\delta_i$ , which is denoted by

$$\delta_1 \gg \dots \gg \delta_s.$$

We define  $\mathcal{P}_{>i} = \{P_{i+1}, \dots, P_s\}$  and

$$\begin{aligned} \Sigma_i &= \{P_i = 0, P_i = \delta_i, P_i = -\delta_i, P_i \geq 2\delta_i, P_i \leq -2\delta_i\}, \\ \Sigma_{\leq i} &= \{\Psi \mid \Psi = \bigwedge_{j=1, \dots, i} \Psi_j, \Psi_j \in \Sigma_j\}. \end{aligned}$$

If  $\Phi$  is a  $(\mathcal{Q}, \mathcal{P})$ -closed formula, we denote by  $\text{Real}_i(\Phi)$  the extension of  $\text{Real}(\Phi)$  to  $\mathbb{R}\langle \delta_1, \dots, \delta_i \rangle^k$ . For  $\Psi \in \Sigma_{\leq i}$ , we denote by  $\text{Real}_i(\Phi \wedge \Psi)$  the intersection of the realization of  $\Psi$  with  $\text{Real}_i(\Phi)$  and by  $b(\Phi \wedge \Psi)$  the sum of the Betti numbers of  $\text{Real}_i(\Phi \wedge \Psi)$ .

**Proposition 7.39.** *For every  $(\mathcal{Q}, \mathcal{P})$ -closed formula  $\Phi$ ,*

$$b(\Phi) \leq \sum_{\substack{\Psi \in \Sigma_{\leq s} \\ \text{Real}_s(\Psi) \subset \text{Real}_s(\Phi)}} b(\Psi).$$

The main ingredient of the proof of the proposition is the following lemma.

**Lemma 7.40.** *For every  $(\mathcal{Q}, \mathcal{P})$ -closed formula  $\Phi$  and every  $\Psi \in \Sigma_{\leq i}$ ,*

$$b(\Phi \wedge \Psi) \leq \sum_{\psi \in \Sigma_{i+1}} b(\Phi \wedge \Psi \wedge \psi).$$

**Proof:** Consider the formulas

$$\begin{aligned} \Phi_1 &= \Phi \wedge \Psi \wedge (P_{i+1}^2 - \delta_{i+1}^2) \geq 0, \\ \Phi_2 &= \Phi \wedge \Psi \wedge (0 \leq P_{i+1}^2 \leq \delta_{i+1}^2). \end{aligned}$$

Clearly,  $\text{Real}_{i+1}(\Phi \wedge \Psi) = \text{Real}_{i+1}(\Phi_1 \vee \Phi_2)$ . Using Proposition 6.44, we have that,

$$b(\Phi \wedge \Psi) \leq b(\Phi_1) + b(\Phi_2) + b(\Phi_1 \wedge \Phi_2).$$

Now, since  $\text{Real}_{i+1}(\Phi_1 \wedge \Phi_2)$  is the disjoint union of

$$\begin{aligned} &\text{Real}_{i+1}(\Phi \wedge \Psi \wedge (P_{i+1} = \delta_{i+1})), \text{Real}_{i+1}(\Phi \wedge \Psi \wedge (P_{i+1} = -\delta_{i+1})), \\ &b(\Phi_1 \wedge \Phi_2) = b(\Phi \wedge \Psi \wedge (P_{i+1} = \delta_{i+1})) + b(\Phi \wedge \Psi \wedge (P_{i+1} = -\delta_{i+1})). \end{aligned}$$

Moreover,

$$\begin{aligned} b(\Phi_1) &= b(\Phi \wedge \Psi \wedge (P_{i+1} \\ &\quad \geq 2\delta_{i+1})) + b(\Phi \wedge \Psi \wedge (P_{i+1} \leq -2\delta_{i+1})), \\ b(\Phi_2) &= b(\Phi \wedge \Psi \wedge (P_{i+1} = 0)). \end{aligned}$$

Indeed, by Theorem 5.46 (Hardt's triviality), denoting

$$F_t = \{x \in \text{Real}_i(\Phi \wedge \Psi) \mid P_{i+1}(x) = t\},$$

there exists  $t_0 \in \mathbb{R}\langle \delta_1, \dots, \delta_i \rangle$  such that

$$F_{[-t_0, 0) \cup (0, t_0]} = \{x \in \text{Reali}_i(\Phi \wedge \Psi) \mid t_0^2 \geq P_{i+1}(x) > 0\}$$

and

$$([-t_0, 0) \times F_{-t_0}) \cup ((0, t_0] \times F_{t_0})$$

are homeomorphic. This implies clearly that

$$F_{[\delta_{i+1}, t_0]} = \{x \in \text{Reali}_{i+1}(\Phi \wedge \Psi) \mid t_0 \geq P_{i+1}(x) \geq \delta_{i+1}\}$$

and

$$F_{[2\delta_{i+1}, t_0]} = \{x \in \text{Reali}_{i+1}(\Phi \wedge \Psi) \mid t_0 \geq P_{i+1}(x) \geq 2\delta_{i+1}\}$$

are homeomorphic, and moreover the homeomorphism can be chosen such that it is the identity on the fibers  $F_{-t_0}$  and  $F_{t_0}$ .

Hence,

$$\mathfrak{b}(\Phi_1) = \mathfrak{b}(\Phi \wedge \Psi \wedge (P_{i+1} \geq 2\delta_{i+1})) + \mathfrak{b}(\Phi \wedge \Psi \wedge (P_{i+1} \leq -2\delta_{i+1})).$$

Note that  $F_0 = \text{Reali}_{i+1}(\Phi \wedge \Psi \wedge (P_{i+1} = 0))$  and  $F_{[-\delta_{i+1}, \delta_{i+1}]} = \text{Reali}_{i+1}(\Phi_2)$ .

Thus, it remains to prove that  $\mathfrak{b}(F_{[-\delta_{i+1}, \delta_{i+1}]}) = \mathfrak{b}(F_0)$ . By Theorem 5.46 (Hardt's triviality), for every  $0 < u < 1$ , there is a fiber preserving semi-algebraic homeomorphism

$$\phi_u: F_{[-\delta_{i+1}, -u\delta_{i+1}]} \rightarrow [-\delta_{i+1}, -u\delta_{i+1}] \times F_{-u\delta_{i+1}}$$

and a semi-algebraic homeomorphism

$$\psi_u: F_{[u\delta_{i+1}, \delta_{i+1}]} \rightarrow [u\delta_{i+1}, \delta_{i+1}] \times F_{u\delta_{i+1}}.$$

We define a continuous semi-algebraic homotopy  $g$  from the identity of  $F_{[-\delta_{i+1}, \delta_{i+1}]}$  to  $\lim_{\delta_{i+1}}$  (from  $F_{[-\delta_{i+1}, \delta_{i+1}]}$  to  $F_0$ ) as follows:

- $g(0, -)$  is  $\lim_{\delta_{i+1}}$ ,
- for  $0 < u \leq 1$ ,  $g(u, -)$  is the identity on  $F_{[-u\delta_{i+1}, u\delta_{i+1}]}$  and sends  $F_{[-\delta_{i+1}, -u\delta_{i+1}]}$  (resp.  $F_{[u\delta_{i+1}, \delta_{i+1}]}$ ) to  $F_{-u\delta_{i+1}}$  (resp.  $F_{u\delta_{i+1}}$ ) by  $\phi_u$  (resp.  $\psi_u$ ) followed by the projection to  $F_{u\delta_{i+1}}$  (resp.  $F_{-u\delta_{i+1}}$ ).

Thus,

$$\mathfrak{b}(F_{[-\delta_{i+1}, \delta_{i+1}]}) = \mathfrak{b}(F_0).$$

Finally,

$$\mathfrak{b}(\Phi \wedge \Psi) \leq \sum_{\psi \in \Sigma_{i+1}} \mathfrak{b}(\Phi \wedge \Psi \wedge \psi). \quad \square$$

**Proof of Proposition 7.39:** Starting from the formula  $\Phi$ , apply Lemma 7.40 with  $\Psi$  the empty formula. Now, repeatedly apply Lemma 7.40 to the terms appearing on the right-hand side of the inequality obtained, noting that for any  $\Psi \in \Sigma_{\leq s}$ ,

- either  $\text{Reali}_s(\Phi \wedge \Psi) = \text{Reali}_s(\Psi)$  and  $\text{Reali}_s(\Psi) \subset \text{Reali}_s(\Phi)$ ,
- or  $\text{Reali}_s(\Phi \wedge \Psi) = \emptyset$ . □

Using an argument analogous to that used in the proof of Theorem 7.30, we prove the following proposition.

**Proposition 7.41.** *For  $0 \leq i \leq k'$ ,*

$$\sum_{\Psi \in \Sigma_{\leq s}} b_i(\Psi) \leq \sum_{j=1}^{k'-i} \binom{s}{j} 6^j d (2d-1)^{k-1}.$$

We first prove the following Lemma 7.42 and Lemma 7.43.

Let  $\mathcal{P} = \{P_1, \dots, P_j\} \subset R[X_1, \dots, X_k]$ , and let  $Q_i = P_i^2(P_i^2 - \delta_i^2)^2(P_i^2 - 4\delta_i^2)$ .  
 Let  $j \geq 1$ ,

$$V'_j = \text{Reali} \left( \bigvee_{1 \leq i \leq j} Q_i = 0, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle) \right),$$

$$W'_j = \text{Reali} \left( \bigvee_{1 \leq i \leq j} Q_i \geq 0, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle) \right).$$

**Lemma 7.42.**

$$b_i(V'_j) \leq (6^j - 1) d (2d - 1)^{k-1}.$$

**Proof:** The set  $\text{Reali}((P_i^2(P_i^2 - \delta_i^2)^2(P_i^2 - 4\delta_i^2) = 0), Z)$  is the disjoint union of

$$\begin{aligned} & \text{Reali}(P_i = 0, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle)), \\ & \text{Reali}(P_i = \delta_i, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle)), \\ & \text{Reali}(P_i = -\delta_i, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle)), \\ & \text{Reali}(P_i = 2\delta_i, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle)), \\ & \text{Reali}(P_i = -2\delta_i, \text{Ext}(Z, R\langle \delta_1, \dots, \delta_j \rangle)). \end{aligned}$$

Moreover, the  $i$ -th Betti number of their union  $V'_j$  is bounded by the sum of the Betti numbers of all possible non-empty sets that can be obtained by taking intersections of these sets using part (a) of Proposition 7.33.

The number of possible  $\ell$ -ary intersection is  $\binom{j}{\ell}$ . Each such intersection is a disjoint union of  $5^\ell$  algebraic sets. The  $i$ -th Betti number of each of these algebraic sets is bounded by  $d(2d-1)^{k-1}$  by Theorem 7.23.

Thus,

$$b_i(V'_j) \leq \sum_{\ell=1}^j \binom{j}{\ell} 5^\ell d (2d - 1)^{k-1} = (6^j - 1) d (2d - 1)^{k-1}. \quad \square$$

**Lemma 7.43.**

$$b_i(W'_j) \leq (6^j - 1) d (2d - 1)^{k-1} + b_i(Z).$$

**Proof:** Let

$$F = \text{Real}\left(\bigwedge_{1 \leq i \leq j} Q_i \leq 0 \vee \bigvee_{1 \leq i \leq j} Q_i = 0, \text{Ext}(Z, \mathbb{R}\langle \delta_1, \dots, \delta_i \rangle)\right).$$

Now,

$$W'_j \cup F = Z, W'_j \cap F = V'_j.$$

Using inequality (6.44) we get that

$$b_i(W'_j) \leq b_i(W'_j \cap F) + b_i(W'_j \cup F) = b_i(V'_j) + b_i(Z)$$

since  $b_i(Z) = b_i(\text{Ext}(Z, \mathbb{R}\langle \delta_1, \dots, \delta_i \rangle))$ . We conclude using Lemma 7.42. □

Now, let

$$S_i = \text{Real}\left(P_i^2(P_i^2 - \delta_i^2)^2(P_i^2 - 4\delta_i^2) \geq 0, \text{Ext}(Z, \mathbb{R}\langle \delta_1, \dots, \delta_s \rangle)\right), 1 \leq i \leq s,$$

and let  $S$  be the intersection of the  $S_i$  with the closed ball in  $\mathbb{R}\langle \delta_1, \dots, \delta_s, \delta \rangle^k$  defined by  $\delta^2 \left( \sum_{1 \leq i \leq k} X_i^2 \right) \leq 1$ . Then, it is clear that

$$\sum_{\Psi \in \Sigma_{\leq s}} b_i(\Psi) = b_i(S).$$

**Proof of Proposition 7.41:** Since, for all  $i < k'$ ,

$$b_i(Z) + b_{k'}(Z) \leq d(2d - 1)^{k-1}$$

by Theorem 7.23 we get that,

$$\sum_{\Psi \in \Sigma_{\leq s}} b_i(\Psi) = b_i(S) \leq \sum_{j=1}^{k'-i} \binom{s}{j} (6^j - 1) d (2d - 1)^{k-1} + \binom{s}{k'-i} b_{k'}(Z)$$

using part (b) of Proposition 7.33 and Lemma 7.43.

Thus, we have that

$$\sum_{\Psi \in \Sigma_{\leq s}} b_i(\Psi) \leq \sum_{j=1}^{k'-i} \binom{s}{j} 6^j d (2d - 1)^{k-1}. \quad \square$$

**Proof of Theorem 7.38:** Theorem 7.38 now follows from Proposition 7.39 and Proposition 7.41. □

### 7.5 Sum of the Betti Numbers of Semi-algebraic Sets

We first describe a construction for replacing any given semi-algebraic subset of a bounded semi-algebraic set by a closed bounded semi-algebraic subset and prove that the new set has the same homotopy type as the original one. Moreover, the polynomials defining the bounded closed semi-algebraic subset are closely related (by infinitesimal perturbations) to the polynomials defining the original subset. In particular, their degrees do not increase, while the number of polynomials used in the definition of the new set is at most twice the square of the number used in the definition of the original set. This construction will be useful later in Chapter 16.

**Definition 7.44.** Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite set of polynomials with  $t$  elements, and let  $S$  be a bounded  $\mathcal{P}$ -closed set. We denote by  $\text{SIGN}(S)$  the set of realizable sign conditions of  $\mathcal{P}$  whose realizations are contained in  $S$ .

Recall that, for  $\sigma \in \text{SIGN}(\mathcal{P})$  we define the level of  $\sigma$  as  $\#\{P \in \mathcal{P} \mid \sigma(P) = 0\}$ . Let,  $\varepsilon_{2t} \gg \varepsilon_{2t-1} \gg \dots \gg \varepsilon_2 \gg \varepsilon_1 > 0$  be infinitesimals, and denote by  $\mathbb{R}_i$  the field  $\mathbb{R}\langle \varepsilon_{2t} \rangle \dots \langle \varepsilon_i \rangle$ . For  $i > 2t$ ,  $\mathbb{R}_i = \mathbb{R}$  and for  $i \leq 0$ ,  $\mathbb{R}_i = \mathbb{R}_1$ .

We now describe the construction. For each level  $m$ ,  $0 \leq m \leq t$ , we denote by  $\text{SIGN}_m(S)$  the subset of  $\text{SIGN}(S)$  of elements of level  $m$ .

Given  $\sigma \in \text{SIGN}_m(\mathcal{P}, S)$ , let  $\text{Reali}(\sigma_+^c)$  be the intersection of  $\text{Ext}(S, \mathbb{R}_{2m})$  with the closed semi-algebraic set defined by the conjunction of the inequalities,

$$\begin{aligned} -\varepsilon_{2m} \leq P \leq \varepsilon_{2m} & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 0, \\ P \geq 0 & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 1, \\ P \leq 0 & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = -1. \end{aligned}$$

and let  $\text{Reali}(\sigma_+^o)$  be the intersection of  $\text{Ext}(S, \mathbb{R}_{2m-1})$  with the open semi-algebraic set defined by the conjunction of the inequalities,

$$\begin{aligned} -\varepsilon_{2m-1} < P < \varepsilon_{2m-1} & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 0, \\ P > 0 & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 1, \\ P < 0 & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = -1. \end{aligned}$$

Notice that, denoting  $\text{Reali}(\sigma)_i = \text{Ext}(\text{Reali}(\sigma), \mathbb{R}_i)$ ,

$$\begin{aligned} \text{Reali}(\sigma)_{2m} & \subset \text{Reali}(\sigma_+^c), \\ \text{Reali}(\sigma)_{2m-1} & \subset \text{Reali}(\sigma_+^o). \end{aligned}$$

Let  $X \subset S$  be a  $\mathcal{P}$ -semi-algebraic set such that

$$X = \bigcup_{\sigma \in \Sigma} \text{Reali}(\sigma)$$

with  $\Sigma \subset \text{SIGN}(S)$ . We denote  $\Sigma_m = \Sigma \cap \text{SIGN}_m(S)$  and define a sequence of sets,  $X^m \subset \mathbb{R}^{tk}$ ,  $0 \leq m \leq t$  inductively by

$$- X^0 = \text{Ext}(X, \mathbb{R}_1).$$

– For  $0 \leq m \leq t$ ,

$$X^{m+1} = \left( X^m \cup \bigcup_{\sigma \in \Sigma_m} \text{Reali}(\sigma_+^c)_1 \right) \setminus \bigcup_{\sigma \in \text{SIGN}_m(S) \setminus \Sigma_m} \text{Reali}(\sigma_+^o)_1,$$

with  $\text{Reali}(\sigma_+^c)_i = \text{Ext}(\text{Reali}(\sigma_+^c), R_i)$ ,  $\text{Reali}(\sigma_+^o)_i = \text{Ext}(\text{Reali}(\sigma_+^o), R_i)$ .

We denote by  $X'$  the set  $X^{t+1}$ . □

**Theorem 7.45.** *The sets  $\text{Ext}(X, R_1)$  and  $X'$  are semi-algebraically homotopy equivalent. In particular,*

$$H_*(X) \cong H_*(X').$$

For the purpose of the proof we introduce several new families of sets defined inductively.

For each  $p$ ,  $0 \leq p \leq t + 1$  we define sets,  $Y_p \subset R_{2p}^k$ ,  $Z_p \subset R_{2p-1}^k$  as follows.

– We define

$$\begin{aligned} Y_p^p &= \text{Ext}(X, R_{2p}) \cup \bigcup_{\sigma \in \Sigma_p} \text{Reali}(\sigma_+^c)_{2p} \\ Z_p^p &= \text{Ext}(Y_p^p, R_{2p-1}) \setminus \bigcup_{\sigma \in \text{SIGN}_p(S) \setminus \Sigma_p} \text{Reali}(\sigma_+^o)_{2p-1}. \end{aligned}$$

– For  $p \leq m \leq t$ , we define

$$\begin{aligned} Y_p^{m+1} &= \left( Y_p^m \cup \bigcup_{\sigma \in \Sigma_m} \text{Reali}(\sigma_+^c)_{2p} \right) \setminus \bigcup_{\sigma \in \text{SIGN}_m(S) \setminus \Sigma_m} \text{Reali}(\sigma_+^o)_{2p} \\ Z_p^{m+1} &= \left( Z_p^m \cup \bigcup_{\sigma \in \Sigma_m} \text{Reali}(\sigma_+^c)_{2p-1} \right) \setminus \bigcup_{\sigma \in \text{SIGN}_m(S) \setminus \Sigma_m} \text{Reali}(\sigma_+^o)_{2p-1}. \end{aligned}$$

We denote by  $Y_p \subset R_{2p}^k$  the set  $Y_p^{t+1}$  and by  $Z_p \subset R_{2p-1}^k$  the set  $Z_p^{t+1}$ .

Note that

- $X = Y_{t+1} = Z_{t+1}$ ,
- $Z_0 = X'$ .

Notice also that for each  $p$ ,  $0 \leq p \leq t$ ,

- $\text{Ext}(Z_{p+1}^{p+1}, R_{2p}) \subset Y_p^p$ ,
- $Z_p^p \subset \text{Ext}(Y_p^p, R_{2p-1})$

The following inclusions follow directly from the definitions of  $Y_p$  and  $Z_p$ .

**Lemma 7.46.** *For each  $p$ ,  $0 \leq p \leq t$ ,*

- $\text{Ext}(Z_{p+1}, R_{2p}) \subset Y_p$ ,
- $Z_p \subset \text{Ext}(Y_p, R_{2p-1})$ .

We now prove that in both the inclusions in Lemma 7.46 above, the pairs of sets are in fact semi-algebraically homotopy equivalent. These suffice to prove Theorem 7.45.

**Lemma 7.47.** *For  $1 \leq p \leq t$ ,  $Y_p$  is semi-algebraically homotopy equivalent to  $\text{Ext}(Z_{p+1}, R_{2p})$ .*

**Proof:** Let  $Y_p(u) \subset R_{2p+1}^k$  denote the set obtained by replacing the infinitesimal  $\varepsilon_{2p}$  in the definition of  $Y_p$  by  $u$ , and for  $u_0 > 0$ , we will denote by

$$Y_p((0, u_0]) = \{(x, u) \mid x \in Y_p(u), u \in (0, u_0]\} \subset R_{2p+1}^{k+1}.$$

By Hardt's triviality theorem there exist  $u_0 \in R_{2p+1}$ ,  $u_0 > 0$  and a homeomorphism,

$$\psi: Y_p(u_0) \times (0, u_0] \rightarrow Y_p((0, u_0]),$$

such that

- $\pi_{k+1}(\phi(x, u)) = u$ ,
- $\psi(x, u_0) = (x, u_0)$  for  $x \in Y_p(u_0)$ ,
- for all  $u \in (0, u_0]$ , and for every sign condition  $\sigma$  on

$$\cup_{P \in \mathcal{P}} \{P, P \pm \varepsilon_{2t}, \dots, P \pm \varepsilon_{2p+1}\},$$

$\psi(\cdot, u)$  defines a homeomorphism of  $\text{Reali}(\sigma, Y_p(u_0))$  to  $\text{Reali}(\sigma, Y_p(u))$ .

Now, we specialize  $u_0$  to  $\varepsilon_{2p}$  and denote the map corresponding to  $\psi$  by  $\phi$ . For  $\sigma \in \Sigma_p$ , we define,  $\text{Reali}(\sigma_{++}^o)$  to be the set defined by

$$\begin{aligned} -2\varepsilon_{2p} < P < 2\varepsilon_{2p} & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 0, \\ P > -\varepsilon_{2p} & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 1, \\ P < \varepsilon_{2p} & \text{ for each } P \in \mathcal{A} \text{ such that } \sigma(P) = -1. \end{aligned}$$

Let  $\lambda: Y_p \rightarrow R_{2p}$  be a semi-algebraic continuous function such that,

$$\begin{aligned} \lambda(x) &= 1 && \text{on } Y_p \cap \cup_{\sigma \in \Sigma_p} \text{Reali}(\sigma_{++}^o), \\ \lambda(x) &= 0 && \text{on } Y_p \setminus \cup_{\sigma \in \Sigma_p} \text{Reali}(\sigma_{++}^o), \\ 0 < \lambda(x) &< 1 && \text{else.} \end{aligned}$$

We now construct a semi-algebraic homotopy,

$$h: Y_p \times [0, \varepsilon_{2p}] \rightarrow Y_p,$$

by defining,

$$\begin{aligned} h(x, t) &= \pi_{1\dots k} \circ \phi(x, \lambda(x)t + (1 - \lambda(x))\varepsilon_{2p}) && \text{for } 0 < t \leq \varepsilon_{2p} \\ h(x, 0) &= \lim_{t \rightarrow 0^+} h(x, t), && \text{else.} \end{aligned}$$

Note that the last limit exists since  $S$  is closed and bounded. We now show that,  $h(x, 0) \in \text{Ext}(Z_{p+1}, R_{2p})$  for all  $x \in Y_p$ .

Let  $x \in Y_p$  and  $y = h(x, 0)$ .

There are two cases to consider.

- $\lambda(x) < 1$ . In this case,  $x \in \text{Ext}(Z_{p+1}, R_{2p})$  and by property (3) of  $\phi$  and the fact that  $\lambda(x) < 1$ ,  $y \in \text{Ext}(Z_{p+1}, R_{2p})$ .



- $\lambda(x) \geq 1$ . Let  $\sigma_y$  be the sign condition of  $\mathcal{P}$  at  $y$  and suppose that  $y \notin \text{Ext}(Z_{p+1}, R_{2p})$ . There are two cases to consider.
  - $\sigma_y \in \Sigma$ . In this case,  $y \in X$  and hence there must exist

$$\tau \in \text{SIGN}_m(S) \setminus \Sigma_m,$$

with  $m > p$  such that  $y \in \text{Reali}(\tau_+^o)$ .

- $\sigma_y \notin \Sigma$ . In this case, taking  $\tau = \sigma_y$ ,  $\text{level}(\tau) > p$  and  $y \in \text{Reali}(\tau_+^o)$ . It follows from the definition of  $y$ , and property (3) of  $\phi$ , that for any  $m > p$ , and every  $\rho \in \text{SIGN}_m(S)$ ,
  - $y \in \text{Reali}(\rho_+^o)$  implies that  $x \in \text{Reali}(\rho_+^o)$ ,
  - $x \in \text{Reali}(\rho_+^c)$  implies that  $y \in \text{Reali}(\rho_+^c)$ .
 Thus,  $x \notin Y_p$  which is a contradiction.

It follows that,

- $h(\cdot, \varepsilon_{2p}): Y_p \rightarrow Y_p$  is the identity map,
- $h(Y_p, 0) = \text{Ext}(Z_{p+1}, R_{2p})$ ,
- $h(\cdot, t)$  restricted to  $\text{Ext}(Z_{p+1}, R_{2p})$  gives a semi-algebraic homotopy between

$$h(\cdot, \varepsilon_{2p})|_{\text{Ext}(Z_{p+1}, R_{2p})} = \text{id}_{\text{Ext}(Z_{p+1}, R_{2p})}$$

and

$$h(\cdot, 0)|_{\text{Ext}(Z_{p+1}, R_{2p})}.$$

Thus,  $Y_p$  is semi-algebraically homotopy equivalent to  $\text{Ext}(Z_{p+1}, R_{2p})$ . □

**Lemma 7.48.** *For each  $p$ ,  $0 \leq p \leq t$ ,  $Z_p$  is semi-algebraically homotopy equivalent to  $\text{Ext}(Y_p, R_{2p-1})$ .*

**Proof:** For the purpose of the proof we define the following new sets for  $u \in R_{2p}$ .

- Let  $Z'_p(u) \subset R_{2p}^k$  be the set obtained by replacing in the definition of  $Z_p$ ,  $\varepsilon_{2j}$  by  $\varepsilon_{2j} - u$  and  $\varepsilon_{2j-1}$  by  $\varepsilon_{2j-1} + u$  for all  $j > p$ , and  $\varepsilon_{2p}$  by  $\varepsilon_{2p} - u$ , and  $\varepsilon_{2p-1}$  by  $u$ . For  $u_0 > 0$  we will denote

$$Z'_p((0, u_0]) = \{(x, u) \mid x \in Z'_p(u), u \in (0, u_0]\}.$$

$Z'_p((0, u_0])$  the set  $\{(x, u) \mid x \in Z'_p(u), u \in (0, u_0]\}$ .

- Let  $Y'_p(u) \subset R_{2p}^k$  be the set obtained by replacing in the definition of  $Y_p$ ,  $\varepsilon_{2j}$  by  $\varepsilon_{2j} - u$  and  $\varepsilon_{2j-1}$  by  $\varepsilon_{2j-1} + u$  for all  $j > p$  and  $\varepsilon_{2p}$  by  $\varepsilon_{2p} - u$ .
- For  $\sigma \in \text{Sign}_m(S)$ , with  $m \geq p$ , let  $\text{Reali}(\sigma_+^c)(u) \subset R_{2p}^k$  denote the set obtained by replacing  $\varepsilon_{2m}$  by  $\varepsilon_{2m} - u$  in the definition of  $\text{Reali}(\sigma_+^c)$ .
- For  $\sigma \in \text{Sign}_m(S)$ , with  $m > p$ , let  $\text{Reali}(\sigma_+^o)(u) \subset R_{2p}^k$  denote the set obtained by replacing  $\varepsilon_{2m-1}$  by  $\varepsilon_{2m-1} + u$  in the definition of  $\text{Reali}(\sigma_+^o)$ .
- Finally, for  $\sigma \in \text{Sign}_p(S)$  let  $\text{Reali}(\sigma_+^o)(u) \subset R_{2p-1}^k$  denote the set obtained by replacing in the definition of  $\text{Reali}(\sigma_+^o)$ ,  $\varepsilon_{2p-1}$  by  $u$ .

Notice that by definition, for any  $u, v \in \mathbb{R}_{2p}$  with  $0 < u \leq v$ ,  $Z'_p(u) \subset Y'_p(u)$ ,  $Z'_p(v) \subset Z'_p(u)$ ,  $Y'_p(v) \subset Y'_p(u)$ , and

$$\bigcup_{0 < s \leq u} Y'_p(s) = \bigcup_{0 < s \leq u} Z'_p(s).$$

We denote by  $Z'_p$  (respectively,  $Y'_p$ ) the set  $Z'_p(\varepsilon_{2p-1})$  (respectively,  $Y'_p(\varepsilon_{2p-1})$ ). It is easy to see that  $Y'_p$  is semi-algebraically homotopy equivalent to  $\text{Ext}(Y_p, \mathbb{R}_{2p-1})$ , and  $Z'_p$  is semi-algebraically homotopy equivalent to  $Z_p$ . We now prove that,  $Y'_p$  is semi-algebraically homotopy equivalent to  $Z'_p$ , which suffices to prove the lemma.

Let  $\mu: Y'_p \rightarrow \mathbb{R}_{2p-1}$  be the semi-algebraic map defined by

$$\mu(x) = \sup_{u \in (0, \varepsilon_{2p-1}]} \{u \mid x \in Z'_p(u)\}.$$

We prove separately (Lemma 7.49 below) that  $\mu$  is continuous. Note that the definition of the set  $Z'_p(u)$  (as well as the set  $Y'_p(u)$ ) is more complicated than the more natural one consisting of just replacing  $\varepsilon_{2p-1}$  in the definition of  $Z_p$  by  $u$ , is due to the fact that with the latter definition the map  $\mu$  defined below is not necessarily continuous.

We now construct a continuous semi-algebraic map,

$$h: Y'_p \times [0, \varepsilon_{2p-1}] \rightarrow Y'_p$$

as follows.

By Hardt's triviality theorem there exist  $u_0 \in \mathbb{R}_{2p}$ , with  $u_0 > 0$  and a semi-algebraic homeomorphism,

$$\psi: Z'_p(u_0) \times (0, u_0] \rightarrow Z'_p((0, u_0]),$$

such that

- $\pi_{k+1}(\psi(x, u)) = u$ ,
- $\psi(x, u_0) = (x, u_0)$  for  $x \in Z'_p(u_0)$ ,
- for all  $u \in (0, u_0]$  and for every sign condition  $\sigma$  of the family,

$$\bigcup_{P \in \mathcal{P}} \{P, P \pm \varepsilon_{2t}, \dots, P \pm \varepsilon_{2p+1}\},$$

the map  $\psi(\cdot, u)$  restricts to a homeomorphism of  $\text{Reali}(\sigma, Z'_p(u_0))$  to  $\text{Reali}(\sigma, Z'_p(u))$ .

We now specialize  $u_0$  to  $\varepsilon_{2p-1}$  and denote by  $\phi$  the corresponding map,

$$\phi: Z'_p \times (0, \varepsilon_{2p-1}] \rightarrow Z'_p((0, \varepsilon_{2p-1}]).$$

Note, that for every  $u$ ,  $0 < u \leq \varepsilon_{2p-1}$ ,  $\phi$  gives a homeomorphism,

$$\phi_u: Z'_p(u) \rightarrow Z'_p.$$

Hence, for every pair,  $u, u', 0 < u \leq u' \leq \varepsilon_{2p-1}$ , we have a homeomorphism,  $\theta_{u,u'}: Z'_p(u) \rightarrow Z'_p(u')$  obtained by composing  $\phi_u$  with  $\phi_{u'}^{-1}$ .

For  $0 \leq u' < u \leq \varepsilon_{2p-1}$ , we let  $\theta_{u,u'}$  be the identity map. It is clear that  $\theta_{u,u'}$  varies continuously with  $u$  and  $u'$ .

For  $x \in Y'_p, t \in [0, \varepsilon_{2p-1}]$  we now define,

$$h(x, t) = \theta_{\mu(x), t}(x).$$

It is easy to verify from the definition of  $h$  and the properties of  $\phi$  listed above that,  $h$  is continuous and satisfies the following.

- $h(\cdot, 0): Y'_p \rightarrow Y'_p$  is the identity map,
- $h(Y'_p, \varepsilon_{2p-1}) = Z'_p$ ,
- $h(\cdot, t)$  restricts to a homeomorphism  $Z'_p \times t \rightarrow Z'_p$  for every  $t \in [0, \varepsilon_{2p-1}]$ .

This proves the required homotopy equivalence. □

We now prove that the function  $\mu$  used in the proof above is continuous.

**Lemma 7.49.** *The semi-algebraic map  $\mu: Y'_p \rightarrow \mathbb{R}_{2p-1}$  defined by*

$$\mu(x) = \sup_{u \in (0, \varepsilon_{2p-1}]} \{u \mid x \in Z'_p(u)\}$$

*is continuous.*

**Proof :** Let  $0 < \delta \ll \varepsilon_{2p-1}$  be a new infinitesimal. In order to prove the continuity of  $\mu$  (which is a semi-algebraic function defined over  $\mathbb{R}_{2p-1}$ ), it suffices, by Proposition 3.5 to show that

$$\lim_{\delta} \text{Ext}(\mu, \mathbb{R}_{2p-1}(\delta))(x') = \lim_{\delta} \text{Ext}(\mu, \mathbb{R}_{2p-1}(\delta))(x)$$

for every pair of points  $x, x' \in \text{Ext}(Y'_p, \mathbb{R}_{2p-1}(\delta))$  such that  $\lim_{\delta} x = \lim_{\delta} x'$ .

Consider such a pair of points  $x, x' \in \text{Ext}(Y'_p, \mathbb{R}_{2p-1}(\delta))$ . Let  $u \in (0, \varepsilon_{2p-1}]$  be such that  $x \in Z'_p(u)$ . We show below that this implies  $x' \in Z'_p(u')$  for some  $u'$  satisfying  $\lim_{\delta} u' = \lim_{\delta} u$ .

Let  $m$  be the largest integer such that there exists  $\sigma \in \Sigma_m$  with  $x \in \text{Reali}(\sigma_+^c)(u)$ . Since  $x \in Z'_p(u)$  such an  $m$  must exist.

We have two cases:

- $m > p$ : Let  $\sigma \in \Sigma_m$  with  $x \in \text{Reali}(\sigma_+^c)(u)$ . Then, by the maximality of  $m$ , we have that for each  $P \in \mathcal{P}$ ,  $\sigma(P) \neq 0$  implies that  $\lim_{\delta} P(x) \neq 0$ . As a result, we have that  $x' \in \text{Reali}(\sigma_+^c)(u')$  for all

$$u' < u - \max_{P \in \mathcal{P}, \sigma(P)=0} |P(x) - P(x')|,$$

and hence we can choose  $u'$  such that  $x' \in \text{Reali}(\sigma_+^c)(u')$  and  $\lim_{\delta} u' = \lim_{\delta} u$ .

- $m \leq p$ : If  $x' \notin Z'_p(u)$  then since  $x' \in Y'_p \subset Y'_p(u)$ ,

$$x' \in \cup_{\sigma \in \text{SIGN}_p(\mathcal{P}, S) \setminus \Sigma_p} \text{Reali}(\sigma_+^o)(u).$$

Let  $\sigma \in \text{SIGN}_p(S) \setminus \Sigma_p$  be such that  $x' \in \text{Reali}(\sigma_+^o)(u)$ . We prove by contradiction that  $\lim_{\delta} \max_{P \in \mathcal{P}, \sigma(P)=0} |P(x')| = u$ .

Assume that

$$\lim_{\delta} \max_{P \in \mathcal{P}, \sigma(P)=0} |P(x')| \neq u.$$

Since,  $x \notin \text{Reali}(\sigma_+^o)(u)$  by assumption, and  $\lim_{\delta} x' = \lim_{\delta} x$ , there must exist  $P \in \mathcal{P}$ ,  $\sigma(P) \neq 0$ , and  $\lim_{\delta} P(x) = 0$ . Letting  $\tau$  denote the sign condition defined by  $\tau(P) = 0$  if  $\lim_{\delta} P(x) = 0$  and  $\tau(P) = \sigma(P)$  else, we have that  $\text{level}(\tau) > p$  and  $x$  belongs to both  $\text{Reali}(\tau_+^o)(u)$  as well as  $\text{Reali}(\tau_+^c)(u)$ .

Now there are two cases to consider depending on whether  $\tau$  is in  $\Sigma$  or not. If  $\tau \in \Sigma$ , then the fact that  $x \in \text{Reali}(\tau_+^c)(u)$  contradicts the choice of  $m$ , since  $m \leq p$  and  $\text{level}(\tau) > p$ . If  $\tau \notin \Sigma$  then  $x$  gets removed at the level of  $\tau$  in the construction of  $Z'_p(u)$ , and hence  $x \in \text{Reali}(\rho_+^c)(u)$  for some  $\rho \in \Sigma$  with  $\text{level}(\rho) > \text{level}(\tau) > p$ . This again contradicts the choice of  $m$ . Thus,  $\lim_{\delta} \max_{P \in \mathcal{P}, \sigma(P)=0} |P(x')| = u$  and since  $x' \notin \cup_{\sigma \in \text{SIGN}_p(\mathcal{C}, S) \setminus \Sigma_p} \text{Reali}(\sigma_+^o)(u')$  for all  $u' < \max_{P \in \mathcal{P}, \sigma(P)=0} |P(x')|$ , we can choose  $u'$  such that  $\lim_{\delta} u' = \lim_{\delta} u$ , and  $x' \notin \cup_{\sigma \in \text{SIGN}_p(\mathcal{P}, S) \setminus \Sigma_p} \text{Reali}(\sigma_+^o)(u')$ .

In both cases we have that  $x' \in Z'_p(u')$  for some  $u'$  satisfying  $\lim_{\delta} u' = \lim_{\delta} u$ , showing that  $\lim_{\delta} \mu(x') \geq \lim_{\delta} \mu(x)$ . The reverse inequality follows by exchanging the roles of  $x$  and  $x'$  in the previous argument. Hence,

$$\lim_{\delta} \mu(x') = \lim_{\delta} \mu(x),$$

proving the continuity of  $\mu$ . □

**Proof of Theorem 7.45:** The theorem follows immediately from Lemmas 7.47 and 7.48. □

We now define the **Betti numbers** of a general  $\mathcal{P}$ -semi-algebraic set and bound them. Given a  $\mathcal{P}$ -semi-algebraic set  $Y \subset \mathbb{R}^k$ , we replace it by

$$X = \text{Ext}(Y, \mathbb{R}(\varepsilon)) \cap \overline{B}_k(0, 1/\varepsilon).$$

Taking  $S = \overline{B}_k(0, 1/\varepsilon)$ , we know by Theorem 7.45 that there is a closed and bounded semi-algebraic set  $X' \subset \mathbb{R}\langle \varepsilon \rangle_1^k$  such that  $\text{Ext}(X, \mathbb{R}\langle \varepsilon \rangle_1)$  and  $X'$  are semi-algebraically homotopy equivalent. We define the Betti numbers  $b_i(Y) := b_i(X')$ . Note that this definition is clearly homotopy invariant since  $Y$  and  $X'$  has are semi-algebraically homotopy equivalent. We denote by  $b(Y) = b(X')$  the sum of the Betti numbers of  $Y$ .

**Theorem 7.50.** *Let  $Y$  be a  $\mathcal{P}$ -semi-algebraic set where  $\mathcal{P}$  is a family of at most  $s$  polynomials of degree  $d$  in  $k$  variables. Then*

$$b(Y) \leq \sum_{i=0}^k \sum_{j=1}^{k-i} \binom{2s^2+1}{j} 6^j d(2d-1)^{k-1}.$$

**Proof:** Take  $S = \overline{B}_k(0, 1/\varepsilon)$  and  $X = \text{Ext}(Y, \mathbb{R}(\varepsilon)) \cap \overline{B}_k(0, 1/\varepsilon)$ . Defining  $X'$  according to Definition 7.44, apply Theorem 7.38 to  $X'$ , noting that the number of polynomials defining  $X'$  is  $2s^2+1$ , and their degrees are bounded by  $d$ .  $\square$

## 7.6 Bibliographical Notes

The inequalities relating the number of critical points of a Morse function (Theorem 7.24) with the Betti numbers of a compact manifold was proved in its full generality by Morse [121], building on prior work by Birkhoff ([25], page 42). Their generalization "with boundary" (Propositions 7.18 and 7.18) can be found in [81, 82]. Using these inequalities, Thom [157], Milnor[118], Oleinik and Petrovsky [124, 125] proved the bound on the sum of the Betti numbers of algebraic sets presented here. Subsequently, using these bounds Warren [165] proved a bound of  $(4esd/k)^k$  on the number of connected components of the realizations of strict sign conditions of a family of polynomials of  $s$  polynomials in  $k$  variables of degree at most  $d$  and Alon [3] derived a bound of  $(8esd/k)^k$  on the number of all realizable sign conditions (not connected components). All these bounds are in terms of the product  $sd$ .

The first result in which the dependence on  $s$  (the combinatorial part) was studied separately from the dependence on  $d$  (algebraic) appeared in [14], where a bound of  $\binom{O(s)}{k'} O(d)^k$  was proved on the number of connected components of all realizable sign conditions of a family of polynomials restricted to variety of dimension  $k'$ . The generalization of the Thom-Milnor bound on the sum of the Betti numbers of basic semi-algebraic sets to more general closed semi-algebraic sets restricted to a variety was first done in [11] and later improved in [17]. The first result bounding the individual Betti numbers of a basic semi-algebraic set appears in [10] and that bounding the individual Betti numbers over all realizable sign conditions in [17]. The construction for replacing any given semi-algebraic subset of a bounded semi-algebraic set by a closed bounded semi-algebraic subset in Section 7.5, as well as the bound on the sum of the Betti numbers of a general semi-algebraic set, appears in [62].

## Complexity of Basic Algorithms

In Section 8.1, we discuss a few notions needed to analyze the complexity of algorithms and illustrate them by several simple examples. In Section 8.2, we study basic algorithms for linear algebra, including computations of determinants and characteristic polynomials of matrices, and signatures of quadratic forms. In Section 8.3, we compute remainder sequences and the related sub-resultant polynomials. The algorithms in this chapter are very basic and will be used throughout the other chapters of the book.

### 8.1 Definition of Complexity

An **algorithm** is a computational procedure that takes an input and after performing a finite number of allowed operations produces an output.

A typical input to an algorithm in this book will be a set of polynomials with coefficients in a ring  $A$  or a matrix with coefficients in  $A$  or a formula involving certain polynomials with coefficients in  $A$ .

Each of our algorithms will depend on a specified **structure**. The specified structure determines which operations are allowed in the algorithm. We list the following structures that will be used most:

- **Ring structure:** the only operations allowed are addition, subtraction, multiplication between elements of a given ring, and deciding whether an element of the ring is zero.
- **Ordered ring structure:** in addition to the ring structure operations, we can also **compare** two elements in a given ordered ring. That is, given  $a, b$  in the ordered ring we can decide whether  $a = b$ ,  $a > b$ , or  $a < b$ .
- **Ring with integer division structure:** in addition to the ring structure operations, it is also possible to do **exact division** by an element of  $\mathbb{Z}$  which can be performed when we know in advance that the result of the division belongs to the ring. In such a ring  $n \cdot 1 \neq 0$  when  $n \in \mathbb{Z}, n \neq 0$ .

- **Integral domain structure:** in addition to the ring structure operations, it is also possible to do **exact division** by an element of a given integral domain which can be performed when we know in advance that the result of a division belongs to the integral domain.
- **Field structure:** in addition to the ring structure operations, we can also perform **division** by any element of a given field, which can be performed only by a non-zero element.
- **Ordered integral domain structure:** in addition to the integral domain structure operations, we can also **compare** two elements of a given ordered integral domain. That is, given  $a, b$  in the ordered integral domain, we can decide whether  $a = b$ ,  $a > b$ , or  $a < b$ ,
- **Ordered field structure:** in addition to the field structure operations, we can also **compare** two elements of a given ordered field.

Which structure is associated to the algorithm will be systematically indicated in the description of the algorithm.

The **size of the input** is always a vector of integers. Typical parameters we use to describe the size of the input are the dimensions of a matrix, the number of polynomials, their degrees, and their number of variables.

The **complexity** of an algorithm in a structure is a function associating to a vector of integers  $v$  describing the size of the input a bound on the number of operations performed by the algorithm in the structure when it runs over all possible inputs of size  $v$ .

*Remark 8.1.* In this definition of complexity, there are many manipulations that are cost free. For example, given a matrix, we can access an element for free. Also the cost of reading the input or writing the output is not taken into account.  $\square$

The same computation has a different complexity depending on the structure which is specified. In a ring  $A$ , the complexity of a single addition or multiplication is 1. However, if the ring  $A$  is  $D[X]$ , where  $D$  is a ring, then the cost of adding two polynomials is one in  $D[X]$ , while the cost of the same operation in  $D$  clearly depends on the degree of the two polynomials.

To illustrate the discussion, we consider first a few basic examples used later in the book.

We consider first arithmetic operations on univariate polynomials.

*Algorithm 8.1.* [Addition of Univariate Polynomials]

- **Structure:** a ring  $A$ .
- **Input:** two univariate polynomials in  $A[X]$ :

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + a_{p-2} X^{p-2} + \cdots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \cdots + b_0. \end{aligned}$$

- **Output:** the sum  $P + Q$ .

- **Complexity:**  $p + 1$ , where  $p$  is a bound on the degree of  $P$  and  $Q$ .
- **Procedure:** For every  $k \leq p$ , compute the coefficient  $c_k$  of  $X^k$  in  $P + Q$ ,

$$c_k := a_k + b_k.$$

Here, the size of the input is one natural number  $p$ , a bound on the degree of the two polynomials. The computation takes place in the ring  $A$ .

*Algorithm 8.2. [Multiplication of Univariate Polynomials]*

- **Structure:** a ring  $A$ .
- **Input:** two univariate polynomials

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + a_{p-2} X^{p-2} + \dots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \dots + b_0. \end{aligned}$$

in  $A[X]$ , with  $p \geq q$ .

- **Output:** the product  $PQ$ .
- **Complexity:**  $(p + 1)(q + 1) + pq$ , where  $p$  is a bound on the degree of  $P$  and  $q$  a bound on the degree of  $Q$ .
- **Procedure:** For each  $k \leq p + q$ , compute the coefficient  $c_k$  of  $X^k$  in  $PQ$ ,

$$c_k := \begin{cases} \sum_{i=0}^k a_{k-i} b_i, & \text{if } 0 \leq k \leq q, \\ \sum_{i=0}^q a_{k-i} b_i, & \text{if } q < k < p, \\ \sum_{i=0}^{k-p} a_{k-i} b_i. & \text{if } p \leq k \leq p + q. \end{cases}$$

Here the size of the input is two natural numbers, a bound on the degree of each of the two polynomials. The computation takes place in the ring  $A$ .

**Complexity analysis:** For every  $k$ ,  $0 \leq k \leq q$ , there are  $k$  additions and  $k + 1$  multiplications in  $A$ , i.e.  $2k + 1$  arithmetic operations. For every  $k$ , such that  $q < k < p$ , there are  $q$  additions and  $q + 1$  multiplications in  $A$ , i.e.  $2q + 1$  arithmetic operations. For every  $k$ ,  $p \leq k < p + q$ , there are  $p + q - k$  additions and  $p + q - k + 1$  multiplications in  $A$ , i.e.  $2(p + q - k) + 1$  arithmetic operations. Since  $\sum_{k=0}^q k = (q + 1)q/2$ ,

$$\sum_{k=0}^q (2k + 1) = \sum_{k=p}^{p+q} (2(p + q - k) + 1) = (q + 1)^2.$$

So there are all together

$$2(q + 1)^2 + (p - q - 1)(2q + 1) = (p + 1)(q + 1) + pq$$

arithmetic operations performed by the algorithm. □

From now on, our estimates on the complexity of an algorithm will often use the notation  $O$ .



**Notation 8.2. [Big O]** Let  $f$  and  $g$  be mappings from  $\mathbb{N}^\ell$  to  $\mathbb{R}$  and  $h$  be a function from  $\mathbb{R}$  to  $\mathbb{R}$ . The expression " $f(v)$  is  $h(O(g(v)))$ " means that there exists a natural number  $b$  such that for all  $v \in \mathbb{N}^\ell$ ,  $f(v) \leq h(bg(v))$ . The expression " $f(v)$  is  $h(\tilde{O}(g(v)))$ " means that there exist natural number  $a$  such that for all  $v \in \mathbb{N}^\ell$ ,  $f(v) \leq h(g(v) \log_2(g(v))^a)$ .  $\square$

For example, the complexity of the algorithms presented for the addition and multiplication of polynomials are  $O(p)$  and  $O(pq)$ .

*Remark 8.3.* The complexity of computing the product of two univariate polynomials depends on the algorithm used. The complexity of the multiplication of two univariate polynomials of degree at most  $d$  is  $O(d^2)$  when the multiplication is done naively, as in Algorithm 8.2,  $O(d^{\log_2(3)})$  when Karatsuba's method is used,  $O(d \log_2(d)) = \tilde{O}(d)$  using the Fast Fourier Transform (FFT). We decided not to enter into these developments and refer the interested reader to [64].  $\square$

*Algorithm 8.3. [Euclidean Division]*

- **Structure:** a field  $K$ .
- **Input:** two univariate polynomials

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + a_{p-2} X^{p-2} + \dots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \dots + b_0. \end{aligned}$$

in  $K[X]$ , with  $b_q \neq 0$ .

- **Output:**  $\text{Quo}(P, Q)$  and  $\text{Rem}(P, Q)$ , the quotient and remainder in the Euclidean division of  $P$  by  $Q$ .
- **Complexity:**  $(p - q + 1)(2q + 3)$ , where  $p$  is a bound on the degree of  $P$  and  $q$  a bound on the degree of  $Q$ .
- **Procedure:**
  - Initialization:  $C := 0$ ,  $R := P$ .
  - For every  $j$  from  $p$  to  $q$ ,

$$\begin{aligned} C &:= C + \frac{\text{cof}_j(R)}{b_q} X^{j-q}, \\ R &:= R - \frac{\text{cof}_j(R)}{b_q} X^{j-q} Q. \end{aligned}$$

- Output  $C, R$ .

Here the size of the input is two natural numbers, a bound on the degree of one polynomial and the degree of the other. The computation takes place in the field  $K$ .

**Complexity analysis:** There are  $p - q + 1$  values of  $j$  to consider. For each value of  $j$ , there is one division,  $q + 1$  multiplications and  $q + 1$  subtractions. Thus, the complexity is bounded by  $(p - q + 1)(2q + 3)$ .  $\square$

The complexity of an algorithm defined in terms of arithmetic operations often does not give a realistic estimate of the actual computation time when the algorithm is implemented. The reason behind this is the intermediate growth of coefficients during the computation. This is why, in the case of integer entries, we also take into account the bitsizes of the integers which occur in the input. The **bitsize** of a non-zero integer is the number of bits in its binary representation. More precisely, the bitsize of  $n$  is  $\tau$  if and only if  $2^{\tau-1} \leq |n| < 2^\tau$ . The bitsize of a rational number is the sum of the bitsizes of its numerators and denominators.

Adding  $n$  integers of bitsizes bounded by  $\tau$  gives an integer of bitsize bounded by  $\tau + \nu$  where  $\nu$  is the bitsize of  $n$ : indeed, if for every  $1 \leq i \leq n$ , we have  $m_i < 2^\tau$ , then  $m_1 + \dots + m_n < n 2^\tau < 2^{\tau+\nu}$ .

Multiplying  $n$  integers of bitsizes bounded by  $\tau$  gives an integer of size bounded by  $n\tau$ : indeed, if for every  $1 \leq i \leq n$ ,  $m_i < 2^\tau$ , then  $m_1 \dots m_n < 2^{n\tau}$ .

When the input of the algorithms belongs to  $\mathbb{Z}$ , it is thus natural to discuss the **binary complexity** of the algorithms, i.e. to estimate the number of bit operations.

Most of the time, the binary complexity of our algorithms is obtained in two steps. First we compute the number of arithmetic operations performed, second we estimate the bitsize of the integers on which these operations are performed. These bitsize estimates do not follow in general from an analysis of the steps of the algorithm itself, but are consequences of bounds coming from the mathematical nature of the objects considered. For example, when all the intermediate results of a computation are determinants of matrices with integer entries, we can make use Hadamard's bound (see Proposition 8.10).

*Remark 8.4.* The binary complexity of an addition of two integers of bitsize  $\tau$  is  $O(\tau)$ . The binary cost of a multiplication of two integers of bitsize  $\tau$  depends strongly of the algorithm used:  $O(\tau^2)$  when the multiplication is done naively,  $O(\tau^{\log_2(3)})$  when Karatsuba's method, is used,  $O(\tau \log_2(\tau) \log_2(\log_2(\tau))) = \tilde{O}(d \tau)$  using FFT. These developments are not included in the book. We refer the interested reader to [64].  $\square$

Now we describe arithmetic operations on multivariate polynomials.

*Algorithm 8.4.* **[Addition of Multivariate Polynomials]**

- **Structure:** a ring  $A$ .
- **Input:** two multivariate polynomials  $P$  and  $Q$  in  $A[X_1, \dots, X_k]$  whose degrees are bounded by  $d$ .
- **Output:** the sum  $P + Q$ .
- **Complexity:**  $\binom{d+k}{k} \leq (d+1)^k$ .
- **Procedure:** For every monomial  $m$  of degree  $\leq d$  in  $k$  variables, denoting by  $a_m$ ,  $b_m$ , and  $c_m$  the coefficients of  $m$  in  $P$ ,  $Q$ , and  $P + Q$ , compute

$$c_m := a_m + b_m.$$

Studying the complexity of this algorithm requires the following lemma.

**Lemma 8.5.** *The number of monomials of degree  $\leq d$  in  $k$  variables is  $\binom{d+k}{k} \leq (d+1)^k$ .*

**Proof:** By induction on  $k$  and  $d$ . The result is true for  $k = 1$  and every  $d$  since there are  $d+1$  monomials of degree less than or equal to  $d$ . Since either a monomial does not depend on  $X_k$  or is a multiple of  $X_k$ , the number of monomials of degree  $\leq d$  in  $k$  variables is the sum of the number of monomials of degree  $\leq d$  in  $k-1$  variables and the number of monomials of degree  $\leq d-1$  in  $k$  variables. Finally, note that  $\binom{d+k-1}{k-1} + \binom{d-1+k}{k} = \binom{d+k}{k}$ .

The estimate  $\binom{d+k}{k} \leq (d+1)^k$  is also proved by induction on  $k$  and  $d$ . The estimate is true for  $k=1$  and every  $d$ , and also for  $d=1$  and every  $k \geq 0$ . Suppose by induction hypothesis that  $\binom{d+k-1}{k-1} \leq (d+1)^{k-1}$  and  $\binom{d-1+k}{k} \leq d^k$ . Then

$$\binom{d+k}{k} \leq (d+1)^{k-1} + d^k \leq (d+1)^{k-1} + d(d+1)^{k-1} = (d+1)^k. \quad \square$$

#### Complexity analysis of Algorithm 8.4:

The complexity is  $\binom{d+k}{k} \leq (d+1)^k$ , using Lemma 8.5, since there is one addition to perform for each  $m$ .

If  $A = \mathbb{Z}$ , and the bitsizes of the coefficients of  $P$  and  $Q$  are bounded by  $\tau$ , the bitsizes of the coefficients of their sum are bounded by  $\tau + 1$ .  $\square$

#### Algorithm 8.5. [Multiplication of Multivariate Polynomials]

- **Structure:** a ring  $A$ .
- **Input:** two multivariate polynomials  $P$  and  $Q$  in  $A[X_1, \dots, X_k]$  whose degrees are bounded by  $p$  and  $q$ .
- **Output:** the product  $PQ$ .
- **Complexity:**  $\leq 2 \binom{p+k}{k} \binom{q+k}{k} \leq 2(p+1)^k (q+1)^k$ .
- **Procedure:** For every monomial  $m$  (resp.  $n$ , resp.  $u$ ) of degree  $\leq p$  (resp.  $\leq q$ , resp.  $\leq p+q$ ) in  $k$  variables, denoting by  $a_m$ ,  $b_n$ , and  $c_u$  the coefficients of  $m$  in  $P$  (resp.  $Q$ , resp.  $P \cdot Q$ ), compute

$$c_u := \sum_{n+m=u} a_n b_m.$$

**Complexity analysis:** Given that there are at most  $\binom{p+k}{k}$  monomials of degree  $\leq p$  and  $\binom{q+k}{k}$  monomials of degree  $\leq q$ , there are at most  $\binom{p+k}{k} \binom{q+k}{k}$  multiplications and  $\binom{p+k}{k} \binom{q+k}{k}$  additions to perform. The complexity is at most  $2 \binom{p+k}{k} \binom{q+k}{k} \leq 2(p+1)^k (q+1)^k$ .

If  $A = \mathbb{Z}$ , and the bitsizes of the coefficients of  $P$  and  $Q$  are bounded by  $\tau$  and  $\sigma$ , the bitsizes of the coefficients of their product are bounded by  $\tau + \sigma + k\nu$  where  $\nu$  is the bitsize of  $p + q + 1$ , since there are at most  $(p+q+1)^k$  monomials of degree  $p+q$  in  $k$  variables.  $\square$

**Algorithm 8.6. [Exact Division of Multivariate Polynomials]**

- **Structure:** a field  $K$ .
- **Input:** two multivariate polynomials  $P$  and  $Q$  in  $K[X_1, \dots, X_k]$  whose degrees are bounded by  $p$  and  $q \leq p$  and such that  $Q$  divides  $P$  in  $K[X_1, \dots, X_k]$ .
- **Output:** the polynomial  $C$  such that  $P = CQ$ .
- **Complexity:**  $\leq \binom{p+k}{k} (2 \binom{q+k}{k} + 1) \leq (2(p+1)^k + 1)(q+1)^k$ .
- **Procedure:**
  - Initialization:  $C := 0, R := P$ .
  - While  $R \neq 0$ , order using the graded lexicographical ordering the monomials of  $P$  and  $Q$  and denote by  $m$  and  $n$  the leading monomial of  $P$  and  $Q$  so obtained. Since  $Q$  divides  $P$ , it is clear that  $n$  divides  $m$ . Denoting by  $a_m$  and  $b_n$  the coefficient of  $m$  and  $n$  in  $P$  and  $Q$ ,

$$C := C + \frac{a_m m}{b_n n}$$

$$R := R - \frac{a_m m}{b_n n} Q.$$

- Output  $C$ .

**Proof of correctness:** The equality  $P = CQ + R$  is maintained throughout the algorithm. Moreover, since  $Q$  divides  $P$ ,  $Q$  divides  $R$ . The algorithm terminates with  $R = 0$ , since the leading monomial of  $R$  decreases strictly for the graded lexicographical ordering in each call to the loop.  $\square$

**Complexity analysis:** There are at most  $\binom{p+k}{k}$  monomials to consider before the loop terminates, and there are for each call to the loop at most one division,  $\binom{q+k}{k}$  multiplications and  $\binom{q+k}{k}$  additions to perform. The complexity is

$$\binom{p+k}{k} \left[ 2 \binom{q+k}{k} + 1 \right] \leq (2(p+1)^k + 1)(q+1)^k.$$

Note that the choice of the leading monomial for the graded lexicographical ordering is cost free in our model of complexity.  $\square$

We consider now how to evaluate a univariate polynomial  $P$  at a value  $b$ .

**Notation 8.6. [Horner]** Let  $P = a_p X^p + \dots + a_0 \in A[X]$ , where  $A$  is a ring. The evaluation process uses the **Horner polynomials associated to  $P$** , which are defined inductively by

$$\begin{aligned} \text{Hor}_0(P, X) &= a_p, \\ &\vdots \\ \text{Hor}_i(P, X) &= X \text{Hor}_{i-1}(P, X) + a_{p-i}. \end{aligned}$$

for  $0 \leq i \leq p$ , so that

$$\text{Hor}_i(P, X) = a_p X^i + a_{p-1} X^{i-1} + \dots + a_{p-i}. \quad (8.1)$$

Note that  $\text{Hor}_p(P, X) = P(X)$ . □

*Algorithm 8.7. [Evaluation of a Univariate Polynomial]*

- **Structure:** a ring  $A$ .
- **Input:**  $P = a_p X^p + \dots + a_0 \in A[X]$  and  $b \in A$ .
- **Output:** the value  $P(b)$ .
- **Complexity:**  $2p$ .
- **Procedure:**
  - Initialize  $\text{Hor}_0(P, b) := a_p$ .
  - For  $i$  from 1 to  $p$ ,

$$\text{Hor}_i(P, b) := b \text{Hor}_{i-1}(P, b) + a_{p-i}.$$

- Output  $\text{Hor}_p(P, b) = P(b)$ .

Here the size of the input is a number, a bound on the degree of  $P$ . The computation takes place in the ring  $A$ .

**Complexity analysis:** The number of arithmetic operations is  $2p$ :  $p$  additions and  $p$  multiplications. □

When the polynomial has coefficients in  $\mathbb{Z}$ , we have the following variant.

*Algorithm 8.8. [Special Evaluation of a Univariate Polynomial]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:**  $P = a_p X^p + \dots + a_0 \in \mathbb{Z}[X]$  and  $b/c \in \mathbb{Q}$  with  $b \in \mathbb{Z}, c \in \mathbb{Z}$ .
- **Output:** the value  $c^p P(b/c)$ .
- **Complexity:**  $4p$ .
- **Procedure:**
  - Initialize  $\bar{H}_0(P, b) := a_p, d := 1$ .
  - For  $i$  from 1 to  $p$ ,

$$\begin{aligned} d &:= cd \\ \bar{\text{Hor}}_i(P, b) &:= b \bar{H}_{i-1}(P, b) + d a_{p-i}. \end{aligned}$$

- Output  $\bar{\text{Hor}}_p(P, b) = c^p P(b/c)$ .

**Complexity analysis:** The number of arithmetic operations is  $4p$ :  $p$  additions and  $3p$  multiplications. If  $\tau$  is a bound on the bitsizes of the coefficients of  $P$  and  $\tau'$  is a bound on the bitsizes of  $b$  and  $c$ , the bitsize of  $\bar{\text{Hor}}_i(P, b)$  is  $\tau + i\tau' + \nu$ , where  $\nu$  is the bitsize of  $p + 1$ , since the bitsize of the product of an integer of bitsize  $\tau$  with  $i$ -times the product of an integer of bitsize  $\tau'$  is  $\tau + i\tau'$ , and the bitsize of the sum of  $i + 1$  numbers of size  $\lambda$  is bounded by  $\lambda + \nu$ . □

The Horner process can also be used for computing the translate of a polynomial.

*Algorithm 8.9. [Translation]*

- **Structure:** a ring  $A$ .
- **Input:**  $P(X) = a_p X^p + \dots + a_0$  in  $A[X]$  and an element  $c \in A$ .
- **Output:** the polynomial  $T = P(X - c)$ .
- **Complexity:**  $p(p+1)$ .
- **Procedure:**
  - Initialization:  $T := a_p$ .
  - For  $i$  from 1 to  $p$ ,
 
$$T := (X - c)T + a_{p-i}.$$
  - Output  $T$ .

**Proof of correctness:** It is immediate to verify that after step  $i$ ,

$$T = a_p (X - c)^i + \dots + a_{p-i}.$$

So after step  $p$ ,  $T = P(X - c)$ . □

**Complexity analysis:** In step  $i$ , the computation of  $(X - c)T$  takes  $i$  multiplications by  $c$  and  $i$  additions (multiplications by  $X$  are not counted). The complexity is the sum of the  $p(p+1)/2$  multiplications by  $c$  and  $p(p+1)/2$  additions and is bounded by  $p(p+1)$ . □

When the polynomial is with coefficients in  $\mathbb{Z}$ , we have the following variant.

*Algorithm 8.10. [Special Translation]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:**  $P(X) = a_p X^p + \dots + a_0$  in  $\mathbb{Z}[X]$  and  $b/c \in \mathbb{Q}$ , with  $b \in \mathbb{Z}$ ,  $c \in \mathbb{Z}$ .
- **Output:** the polynomial  $c^p P(X - b/c)$ .
- **Complexity:**  $3p(p+3)/2$ .
- **Procedure:**
  - Initialization:  $\bar{T}_0 := a_p$ ,  $d := 1$ .
  - For  $i$  from 1 to  $p$ ,
 
$$d := cd$$

$$\bar{T}_i := (cX - b)\bar{T}_{i-1} + d \cdot a_{p-i}.$$
  - Output  $\bar{T}_p$ .

**Proof of correctness:** It is immediate to verify that after step  $i$ ,

$$\bar{T}_i = c^i (a_p (X - b/c)^i + \dots + a_{p-i}).$$

So after step  $p$ ,  $\bar{T}_p = c^p P(X - b/c)$ . □

**Complexity analysis:** In step  $i$ , the computation of  $\bar{T}$  takes  $2i + 2$  multiplications and  $i$  additions. The complexity is the sum of the  $p(p+3)$  multiplications and  $p(p+1)/2$  additions and is bounded by  $3p(p+3)/2$ .

Let  $\tau$  be a bound on the bitsizes of the coefficients of  $P$ ,  $\tau'$  a bound on the bitsizes of  $b$  and  $c$ , and  $\nu$  is the bitsize of  $p + 1$ . Since

$$\sum_{k=0}^i a_{p-k} (bX - c)^{i-k} = \sum_{k=0}^i a_{p-k} \binom{j}{i-k} b^j (-c)^{i-k-j} X^j,$$

the bitsizes of the coefficients of  $\bar{T}_i$  is  $\tau + i(1 + \tau') + \nu$ : the bitsize of a binomial coefficient  $\binom{i-k}{j}$  is at most  $i$ , the bitsize of the product of an integer of bitsize  $\tau$  with the product of  $i - k$  integers of bitsize  $\tau'$  is bounded by  $\tau + i\tau'$ , and the bitsize of the sum of  $i + 1$  numbers of size  $\lambda$  is bounded by  $\lambda + \nu$ .  $\square$

*Remark 8.7.* Using fast arithmetic, a translation by 1 in a polynomial of degree  $d$  and bit size  $\tau$  can be computed with binary complexity  $\tilde{O}(d\tau)$  [64].  $\square$

We give an algorithm computing the coefficients of a polynomial knowing its Newton sums.

*Algorithm 8.11. [Newton Sums]*

- **Structure:** a ring  $D$  with division in  $\mathbb{Z}$ .
- **Input:** the Newton sums  $N_i$ ,  $i = 0, \dots, p$ , of a monic polynomial

$$P = X^p + a_{p-1}X^{p-1} + \dots + a_0$$

in  $D[X]$ .

- **Output:** the list of coefficients  $1, a_{p-1}, \dots, a_0$  of  $P$ .
- **Complexity:**  $p(p + 1)$ .
- **Procedure:**
  - $a_p := 1$ .
  - For  $i$  from 1 to  $p$ ,

$$a_{p-i} := \frac{-1}{i} \left( \sum_{j=1}^i a_{p-i+j} N_j \right).$$

**Proof of correctness:** Follows from Equation (4.1). Note that we have to know in advance that  $P \in D[X]$ .  $\square$

**Complexity analysis:** The computation of each  $a_{p-i}$  takes  $2i + 1$  arithmetic operations in  $D$ . Since the complexity is bounded by

$$\sum_{i=1}^p (2i + 1) = 2 \frac{p(p-1)}{2} + p = p(p + 1). \quad \square$$

Note also that the Newton formulas (Equation (4.1)) could also be used to compute the Newton sums from the coefficients.

We end this list of examples with arithmetic operations on matrices.

*Algorithm 8.12. [Addition of Matrices]*

- **Structure:** a ring  $A$ .

- **Input:** two  $n \times m$  matrices  $M = [m_{i,j}]$  and  $N = [n_{i,j}]$  with entries in  $A$ .
- **Output:** the sum  $S = [s_{i,j}]$  of  $M$  and  $N$ .
- **Complexity:**  $nm$ .
- **Procedure:** For every  $i, j$ ,  $i \leq n$ ,  $j \leq m$ ,

$$s_{i,j} := m_{i,j} + n_{i,j}.$$

Here the size of the input is two natural numbers  $n, m$ . The computation takes place in the ring  $A$ .

**Complexity analysis:** The complexity is  $nm$  in  $A$  since there are  $nm$  entries to compute and each of them is computed by one single addition.

If  $A = \mathbb{Z}$ , and the bitsizes of the entries of  $M$  and  $N$  are bounded by  $\tau$ , the bitsizes of the entries of their sum are bounded by  $\tau + 1$ .

If  $A = \mathbb{Z}[Y]$ ,  $Y = Y_1, \dots, Y_t$ , and the degrees in  $Y$  of the entries of  $M$  and  $N$  are bounded by  $c$ , while the bitsizes of the entries of  $M$  and  $N$  are bounded by  $\tau$ , the degrees in  $Y$  of the entries of their sum is bounded by  $c$ , and the bitsizes of the coefficients of the entries of their sum are bounded by  $\tau + 1$ .  $\square$

*Algorithm 8.13. [Multiplication of Matrices]*

- **Structure:** a ring  $A$ .
- **Input:** two matrices  $M = [m_{i,j}]$  and  $N = [n_{j,k}]$  of size  $n \times m$  and  $m \times \ell$  with entries in  $A$ .
- **Output:** the product  $P = [p_{i,k}]$  of  $M$  and  $N$ .
- **Complexity:**  $n\ell(2m - 1)$ .
- **Procedure:** For each  $i, k$ ,  $i \leq n$ ,  $k \leq \ell$ ,

$$p_{i,k} = \sum_{j=1}^m m_{i,j} n_{j,k}.$$

**Complexity analysis:** For each  $i, k$  there are  $m$  multiplications and  $m - 1$  additions. The complexity is  $n\ell(2m - 1)$ .

If  $A = \mathbb{Z}$ , and the bitsizes of the entries of  $M$  and  $N$  are bounded by  $\tau$  and  $\sigma$ , the bitsizes of the entries of their product are bounded by  $\tau + \sigma + \mu$ , where  $\mu$  is the bitsize of  $m$ .

If  $A = \mathbb{Z}[Y]$ ,  $Y = Y_1, \dots, Y_k$ , and the degrees in  $Y$  of the entries of  $M$  and  $N$  are bounded by  $p$  and  $q$ , while the bitsizes of the entries of  $M$  and  $N$  are bounded by  $\tau$  and  $\sigma$ , the degrees in  $Y$  of the entries of their product are bounded by  $p + q$ , and the bitsizes of the coefficients of the entries of their product are bounded by  $\tau + \sigma + k\nu + \mu$  where  $\mu$  is the bitsize of  $m$  and  $\nu$  is the bitsize of  $p + q + 1$ , since the number of monomials of degree  $p + q$  in  $k$  variables is bounded by  $(p + q + 1)^k$ .  $\square$



**Algorithm 8.14. [Multiplication of Several Matrices]**

- **Structure:** a ring  $A$ .
- **Input:**  $m$  matrices  $M_1 \dots M_m$  of size  $n \times n$ , with entries in  $A$ .
- **Output:** the product  $P$  of  $M_1, \dots, M_m$ .
- **Complexity:**  $(m-1)n^2(2n-1)$ .
- **Procedure:** Initialize  $N_1 := M_1$ . For  $i$  from 2 to  $m$  define  $N_i = N_{i-1} M_i$ .

**Complexity analysis:** For each  $i$  from 2 to  $m$ , and  $j, k$  from 1 to  $n$ , there are  $n$  multiplications and  $n-1$  additions. The complexity is  $(m-1)n^2(n-1)$ .

If  $A = \mathbb{Z}$ , and the bitsizes of the entries of the  $M_i$  are bounded by  $\tau$ , the bitsizes of the entries of their product are bounded by  $m(\tau + \mu)$  where  $\mu$  is the bitsize of  $n$ .

If  $A = \mathbb{Z}[Y]$ ,  $Y = Y_1, \dots, Y_k$ , and the degrees in  $Y$  of the entries of the  $M_i$  are bounded by  $p$ , while the bitsizes of the entries of the  $M_i$  are bounded by  $\tau$ , the degrees in  $Y$  of the entries of their product are bounded by  $mp$ , and the bitsizes of the coefficients of the entries of their product are bounded by  $m(\tau + \mu) + k\nu$  where  $\mu$  is the bitsize of  $n$  and  $\nu$  is the bitsize of  $kp+1$ .  $\square$

*Remark 8.8.* The complexity of computing the product of two matrices depends on the algorithm used. The complexity of the multiplication of two square matrices of size  $n$  is  $O(n^3)$  when the multiplication is done naively, as in Algorithm 8.13,  $O(n^{\log_2(7)})$  when Strassen's method is used. Even more efficient algorithms are known but we have decided not to include this topic in this book. The interested reader is referred to [64].

Similar remarks were made earlier for the multiplications of polynomials and of integers, and apply also to the euclidean remainder sequence and to most of the algorithms dealing with univariate polynomials and linear algebra presented in Chapters 8 and 9. Explaining sophisticated algorithms would have required a lot of effort and many more pages. In order to prove the complexity estimates we present in Chapters 10 to 15, complexities of  $n^{O(1)}$  for algorithms concerning univariate polynomials and linear algebra (where  $n$  is a bound on the degrees or on the size of the matrices) are sufficient.  $\square$

## 8.2 Linear Algebra

### 8.2.1 Size of Determinants

**Proposition 8.9. [Hadamard]** *Let  $M$  be an  $n \times n$  matrix with integer entries. Then the determinant of  $M$  is bounded by the product of the euclidean norms of the columns of  $M$ .*

**Proof:** If  $\det(M) = 0$ , the result is certainly true. Otherwise, the column vectors of  $M$ ,  $v_1, \dots, v_n$ , span  $\mathbb{R}^n$ . We denote by  $u \cdot v$  the inner product of  $u$  and  $v$ . Using the Gram-Schmidt orthogonalization process (Proposition 4.40), there are vectors  $w_1, \dots, w_n$  with the following properties

- $w_i - v_i$  belong to the vector space spanned by  $w_1, \dots, w_{i-1}$ ,
- $\forall i \forall j \ j \neq i, w_i \cdot w_j = 0$ .

Moreover, denoting  $u_i = w_i - v_i$ ,

$$\begin{aligned} \|w_i\|^2 + \|u_i\|^2 &= \|v_i\|^2, \\ \|w_i\| &\leq \|v_i\|. \end{aligned}$$

Then it is clear that

$$|\det(M)| = \prod_{i=1}^n \|w_i\| \leq \prod_{i=1}^n \|v_i\|. \quad \square$$

**Corollary 8.10.** *Let  $M$  be an  $n \times n$  matrix with integer entries of bitsizes at most  $\tau$ . Then the bitsize of the determinant of  $M$  is bounded by  $n(\tau + \nu/2)$ , where  $\nu$  is the bitsize of  $n$ .*

**Proof:** If  $n < 2^\nu$  and  $|m_{i,j}| < 2^\tau$  then  $\sqrt{\sum_{i=1}^n m_{i,j}^2} < \sqrt{n}2^\tau < 2^{\tau+\nu/2}$ .

Thus  $|\det(M)| < 2^{n(\tau+\nu/2)}$ , using Lemma 8.9. □

The same kind of behavior is observed when we consider degrees of polynomials rather than bitsize. Things are even simpler, since there is no carry to take into account in the degree estimates.

**Proposition 8.11.** *Let  $M$  be an  $n \times n$  matrix with entries that are polynomials in  $Y_1, \dots, Y_k$  of degrees  $d$ . Then the determinant considered as a polynomial in  $Y_1, \dots, Y_k$  has degree in  $Y_1, \dots, Y_k$  bounded by  $dn$ .*

**Proof:** This follows from  $\det(M) = \sum_{\sigma \in \mathcal{S}_n} (-1)^{\varepsilon(\sigma)} \prod_{i=1}^n m_{\sigma(i),i}$ , where  $\varepsilon(\sigma)$  is the signature of  $\sigma$ . □

Moreover we have

**Proposition 8.12.** *Let  $M$  be an  $n \times n$  matrix with entries that are polynomials in  $Y_1, \dots, Y_k$  of degrees  $d$  in  $Y_1, \dots, Y_k$  and coefficients in  $\mathbb{Z}$  of bitsize  $\tau$ . Then the determinant considered as a polynomial in  $Y_1, \dots, Y_k$  has degrees in  $Y_1, \dots, Y_k$  bounded by  $dn$ , and coefficients of bitsize  $(\tau + \nu)n + k\mu$  where  $\nu$  is the bitsize of  $n$  and  $\mu$  is the bitsize of  $nd + 1$ .*

**Proof:** The only thing which remains to prove is the result on the bitsize. Performing the multiplication of  $n$  monomials appearing in the entries of the matrix produces integers of bitsize  $\tau n$ . Since the number of monomials of a polynomial of degree  $nd$  in  $k$  variables is bounded by  $(nd + 1)^k$  by Lemma 8.5, the bitsizes of the coefficients of the products of  $n$  entries of the matrix are bounded by  $(\tau + \nu)n + k\mu$ . Since there are  $n!$  terms in the determinant, and the bitsize of  $n!$  is bounded by  $n\nu$  the final bound is  $(\tau + \nu)n + k\mu$ . □

### 8.2.2 Evaluation of Determinants

The following method, which is the standard row reduction technique, can be used to compute the determinant of a square matrix with coefficients in a field.

*Algorithm 8.15. [Gauss]*

- **Structure:** a field  $K$ .
- **Input:** an  $n \times n$  matrix  $M = [m_{i,j}]$  with coefficients in  $K$ .
- **Output:** the determinant of  $M$ .
- **Complexity:**  $O(n^3)$ .
- **Procedure:**
  - Initialization:  $k := 0$  and  $g_{i,j}^{(0)} := m_{i,j}$ .
  - For  $k$  from 0 to  $n - 2$ ,
    - If for every  $j = k + 1, \dots, n$ ,  $g_{k+1,j}^{(k)} = 0$ , output  $\det(M) = 0$ .
    - Otherwise, exchanging columns if needed, suppose  $g_{k+1,k+1}^{(k)} \neq 0$ .
    - For  $i$  from  $k + 2$  to  $n$ ,
 
$$g_{i,k+1}^{(k+1)} := 0,$$
    - For  $j$  from  $k + 2$  to  $n$ ,

$$g_{i,j}^{(k+1)} := g_{i,j}^{(k)} - \frac{g_{i,k+1}^{(k)}}{g_{k+1,k+1}^{(k)}} g_{k+1,j}^{(k)}. \quad (8.2)$$

- Output

$$\det(M) = (-1)^s g_{1,1}^{(0)} \cdots g_{n,n}^{(n-1)} \quad (8.3)$$

(where  $s$  is the number of exchanges of columns in the intermediate computations).

*Example 8.13.* Consider the following matrix

$$M := \begin{bmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{bmatrix},$$

and suppose  $a_1 \neq 0$  and  $b_2 a_1 - b_1 a_2 \neq 0$ . Performing the first step of Algorithm 8.15 (Gauss), we get

$$\begin{aligned} g_{22}^{(1)} &= \frac{a_1 b_2 - b_1 a_2}{a_1} \\ g_{23}^{(1)} &= \frac{a_1 c_2 - c_1 a_2}{a_1} \\ g_{32}^{(1)} &= \frac{a_1 b_3 - b_1 a_3}{a_1} \\ g_{33}^{(1)} &= \frac{a_1 c_3 - c_1 a_3}{a_1} \end{aligned}$$

After the first step of reduction we have obtained the matrix

$$M_1 = \begin{bmatrix} a_1 & b_1 & c_1 \\ 0 & g_{22}^{(1)} & g_{23}^{(1)} \\ 0 & g_{32}^{(1)} & g_{33}^{(1)} \end{bmatrix}.$$

Note that the determinant of  $M_1$  is the same as the determinant of  $M$  since  $M_1$  is obtained from  $M$  by adding a multiple of the first row to the second and third row.

Performing the second step of Algorithm 8.15 (Gauss), we get

$$g_{33}^{(2)} = \frac{c_3 a_1 b_2 - c_3 b_1 a_2 - c_1 a_3 b_2 - c_2 a_1 b_3 + c_2 b_1 a_3 + c_1 a_2 b_3}{b_2 a_1 - b_1 a_2}$$

After the second step of reduction we have obtained the triangular matrix

$$M' = \begin{bmatrix} a_1 & b_1 & c_1 \\ 0 & g_{22}^{(1)} & g_{23}^{(1)} \\ 0 & 0 & g_{33}^{(2)} \end{bmatrix}.$$

Note that the determinant of  $M'$  is the same as the determinant of  $M$  since  $M'$  is obtained from  $M_1$  by adding a multiple of the second row to the third row.

Finally, since  $g_{11}^{(0)} = a_1$ ,

$$\det(M) = \det(M') = g_{11}^{(0)} g_{22}^{(1)} g_{33}^{(2)}. \quad \square$$

**Proof of correctness:** The determinant of the  $n \times n$  matrix  $M' = [g_{i,j}^{(i-1)}]$  obtained at the end of the algorithm is equal to the determinant of  $M$  since the determinant does not change when a multiple of another row is added to a row. Thus, taking into account exchanges of rows,

$$\det(M) = \det(M') = (-1)^s g_{11}^{(0)} \cdots g_{nn}^{(n-1)}. \quad \square$$

**Complexity analysis:** The number of calls to the main loop are at most  $n - 1$ , the number of elements computed in each call to the loop is at most  $(n - i)^2$ , and the computation of an element is done by 3 arithmetic operations. So, the complexity is bounded by

$$3 \left( \sum_{i=1}^{n-1} (n-i)^2 \right) = \frac{2n^3 - 3n^2 + n}{2} = O(n^3).$$

Note that if we are interested only in the bound  $O(n^3)$ , we can estimate the number of elements computed in each call to the loop by  $n^2$  since being more precise changes the constant in front of  $n^3$  but not the fact that the complexity is bounded by  $O(n^3)$ .  $\square$

*Remark 8.14.* It is possible, using other methods, to compute determinants of matrices of size  $n$  in parallel complexity  $O(\log_2(n)^2)$  using  $n^{O(1)}$  processors [127]. As a consequence it is possible to compute them in complexity  $n^{O(1)}$ , using only  $\log_2(n)^{O(1)}$  space at a given time [127].  $\square$

As we can see in Example 8.13, it is annoying to see denominators arising in a determinant computation, since the determinant belongs to the ring generated by the entries of the matrix. This is fixed in what follows.

**Notation 8.15. [Bareiss]** Let  $M_{i,j}^{(k)}$  be the  $(k+1) \times (k+1)$  matrix obtained by taking

$$\begin{cases} m_{i',j'}^{(k)} = m_{i',j'} & \text{for } i' = 1, \dots, k, j' = 1, \dots, k, \\ m_{k+1,j'}^{(k)} = m_{i,j'} & \text{for } j' = 1, \dots, k, \\ m_{i',k+1}^{(k)} = m_{i',j} & \text{for } i' = 1, \dots, k, \\ m_{k+1,k+1}^{(k)} = m_{i,j}. \end{cases}$$

and define  $b_{i,j}^{(k)} = \det(M_{i,j}^{(k)})$ . Then  $b_{k,k}^{(k-1)}$  is the **principal  $k$ -th minor** of  $M$ , i.e. the determinant of the submatrix extracted from  $M$  on the  $k$  first rows and columns. It follows from the definition of the  $b_{i,j}^{(k)}$  that if  $M$  has entries in an integral domain  $D$  then  $b_{i,j}^{(k)} \in D$ .  $\square$

In the following discussion, we always suppose without loss of generality that if  $b_{k+1,k+1}^{(k)} = 0$  then  $b_{k+1,j}^{(k)} = 0$  for  $j = k+2, \dots, n$ , since this condition is fulfilled after a permutation of columns.

Note that by (8.3), if  $i, j \geq k+1$ ,

$$b_{i,j}^{(k)} = g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)} g_{i,j}^{(k)}. \tag{8.4}$$

Indeed, denoting by  $g'_{i,j}^{(k)}$  the output of Gauss's method applied to  $M_{i,j}^{(k)}$ , it is easy to check that  $g'_{i,i}^{(i-1)} = g_{i,i}^{(i-1)}$  for  $i = 1, \dots, k$ , and  $g'_{k+1,k+1}^{(k)} = g_{i,j}^{(k)}$ .

**Proposition 8.16.**

$$b_{i,j}^{(k+1)} = \frac{b_{k+1,k+1}^{(k)} b_{i,j}^{(k)} - b_{i,k+1}^{(k)} b_{k+1,j}^{(k)}}{b_{k,k}^{(k-1)}}.$$

**Proof:** The result follows easily from the recurrence (8.2) and equation (8.4). Indeed (8.4) implies

$$\begin{aligned} \frac{b_{k+1,k+1}^{(k)} b_{i,j}^{(k)} - b_{i,k+1}^{(k)} b_{k+1,j}^{(k)}}{b_{k,k}^{(k-1)}} &= \frac{(g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)})^2 (g_{k+1,k+1}^{(k)} g_{i,j}^{(k)} - g_{i,k+1}^{(k)} g_{k+1,j}^{(k)})}{g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)}} \\ &= g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)} (g_{k+1,k+1}^{(k)} g_{i,j}^{(k)} - g_{i,k+1}^{(k)} g_{k+1,j}^{(k)}). \end{aligned}$$

On the other hand, (8.2) implies that

$$g_{k+1,k+1}^{(k)} g_{i,j}^{(k)} - g_{i,k+1}^{(k)} g_{k+1,j}^{(k)} = g_{k+1,k+1}^{(k)} g_{i,j}^{(k+1)}. \tag{8.5}$$

So

$$\frac{b_{k+1,k+1}^{(k)} b_{i,j}^{(k)} - b_{i,k+1}^{(k)} b_{k+1,j}^{(k)}}{b_{k,k}^{(k-1)}} = g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)} g_{k+1,k+1}^{(k)} g_{i,j}^{(k+1)}.$$

Using again (8.4),

$$g_{1,1}^{(0)} \cdots g_{k,k}^{(k-1)} g_{k+1,k+1}^{(k)} g_{i,j}^{(k+1)} = b_{i,j}^{(k+1)}, \tag{8.6}$$

and the result follows. □

Note that (8.6) implies that, if  $b_{k+1,k+1}^{(k)} \neq 0$ ,

$$g_{i,j}^{(k+1)} = \frac{b_{i,j}^{(k+1)}}{b_{k+1,k+1}^{(k)}}. \tag{8.7}$$

A new algorithm for computing the determinant follows from Proposition 8.16.

*Algorithm 8.16. [Dogdson-Jordan-Bareiss]*

- **Structure:** a domain  $D$ .
- **Input:** an  $n \times n$  matrix  $M = [m_{i,j}]$  with coefficients in  $D$ .
- **Output:** the determinant of  $M$ .
- **Complexity:**  $O(n^3)$ .
- **Procedure:**
  - Initialization:  $k := 0$  and  $b_{i,j}^{(0)} := m_{i,j}$ ,  $b_{0,0}^{(-1)} := 1$ .
  - For  $k$  from 0 to  $n - 2$ ,
    - If for every  $j = k + 1, \dots, n$ ,  $b_{k+1,j}^{(k)} = 0$ , output  $\det(M) = 0$ .
    - Otherwise, exchanging columns if needed, suppose that  $b_{k+1,k+1}^{(k)} \neq 0$ .
    - For  $i$  from  $k + 2$  to  $n$ ,
 
$$b_{i,k+1}^{(k+1)} := 0,$$
    - For  $j$  from  $k + 2$  to  $n$ ,

$$b_{i,j}^{(k+1)} := \frac{b_{k+1,k+1}^{(k)} b_{i,j}^{(k)} - b_{i,k+1}^{(k)} b_{k+1,j}^{(k)}}{b_{k,k}^{(k-1)}}. \tag{8.8}$$

- Output

$$\det(M) = (-1)^s b_{n,n}^{(n-1)} \tag{8.9}$$

(where  $s$  is the number of exchanges of columns in the intermediate computation)

**Proof of correctness:** The correctness follows from Proposition 8.16. Note that although divisions are performed, they are always exact divisions, since we know from Proposition 8.16 that all the intermediate computations obtained by a division in the algorithm are determinants extracted from  $M$  and hence belong to  $D$ .  $\square$

**Complexity analysis:** The number of calls to the main loop are at most  $n - 1$ , the number of elements computed in each call to the loop is at most  $(n - i)^2$ , and the computation of an element is done by 4 arithmetic operations. So, the complexity is bounded by

$$4 \left( \sum_{i=1}^{n-1} (n-i)^2 \right) = \frac{4n^3 - 6n^2 + 2n}{3} = O(n^3).$$

If  $M$  is a matrix with integer coefficients having bitsize at most  $\tau$ , the arithmetic operations in the algorithm are performed on integers of bit-size  $n(\tau + \nu)$ , where  $\nu$  is the bitsize of  $n$ , using Hadamard's bound (Corollary 8.10).  $\square$

*Example 8.17.* Consider again

$$M := \begin{bmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{bmatrix}.$$

Performing the first step of Algorithm 8.16 (Dogdson-Jordan-Bareiss), we get

$$\begin{aligned} b_{22}^{(1)} &= a_1 b_2 - b_1 a_2, \\ b_{23}^{(1)} &= a_1 c_2 - c_1 a_2, \\ b_{32}^{(1)} &= a_1 b_3 - b_1 a_3, \\ b_{33}^{(1)} &= a_1 c_3 - c_1 a_3. \end{aligned}$$

which are determinants extracted from  $M$ .

Performing the second step of Algorithm 8.16 (Dogdson-Jordan-Bareiss), we get

$$\begin{aligned} b_{33}^{(2)} &= \frac{(a_1 b_2 - b_1 a_2)(a_1 c_3 - c_1 a_3) - (a_1 c_2 - c_1 a_2)(a_1 b_3 - b_1 a_3)}{a_1} \\ &= c_3 a_1 b_2 - c_3 b_1 a_2 - c_1 a_3 b_2 - c_2 a_1 b_3 + c_2 b_1 a_3 + c_1 a_2 b_3. \end{aligned}$$

Finally,

$$\det(M) = b_{33}^{(2)}. \quad \square$$

*Remark 8.18.* It is easy to see than either Algorithm 8.15 (Gauss) or Algorithm 8.16 (Dogdson-Jordan-Bareiss) can be adapted as well to compute the rank of the matrix with the same complexity.  $\square$

**Exercise 8.1.** Describe algorithms for computing the rank of a matrix by adapting Algorithm 8.15 (Gauss) and Algorithm 8.16 (Dogdson-Jordan-Bareiss).

### 8.2.3 Characteristic Polynomial

Let  $A$  be a ring and  $M$  be a matrix  $M = (m_{ij}) \in A^{n \times n}$ . The first idea of the method we present to compute the characteristic polynomial is to compute the traces of the powers of  $M$ , and to use Algorithm 8.11 (Newton Sums) to recover the characteristic polynomial. Indeed the trace of  $M^i$  is the  $i$ -th Newton sum of the characteristic polynomial of  $M$ . The second idea is to notice that, in order to compute the trace of a product of two matrices  $M$  and  $N$ , it is not necessary to compute the product  $MN$ , since  $\text{Tr}(MN) = \sum_{k,\ell} m_{k,\ell} n_{\ell,k}$ . So rather than computing all the powers  $M^i$  of  $M$ ,  $i = 2, \dots, n$  then all the corresponding traces, it is enough, defining  $r$  as the smallest integer  $> \sqrt{n}$ , to compute the powers  $M^i$  for  $i = 2, \dots, r - 1$ , the powers  $M^{jr}$  for  $j = 2, \dots, r - 1$  and then  $\text{Tr}(M^{rj+i}) = \text{Tr}(M^i M^{jr})$ .

*Algorithm 8.17. [Characteristic Polynomial]*

- **Structure:** a ring with integer division  $A$ .
- **Input:** an  $n \times n$  matrix  $M = [m_{i,j}]$ , with coefficients in  $A$ .
- **Output:**  $\text{CharPol}(M) = \det(X \text{Id}_n - M)$ , the characteristic polynomial of  $M$ .
- **Complexity:**  $O(n^{3.5})$ .
- **Procedure:**
  - Define  $r$  as the smallest integer  $> \sqrt{n}$ .
  - Computation of powers  $M^i$  for  $i < r$  and their traces.
    - $B_0 := \text{Id}_n, N_0 := n$ .
    - For  $i$  from 0 to  $r - 2$ 

$$B_{i+1} := MB_i, N_{i+1} := \text{Tr}(B_{i+1}).$$
  - Computation of powers  $M^{rj}$  for  $j < r$  and their traces.
    - $C_1 := MB_{r-1}, N_r = \text{Tr}(C_1)$ .
    - For  $j$  from 1 to  $r - 2$ 

$$C_{j+1} = C_1 C_j, N_{(j+1)r} = \text{Tr}(C_{j+1}).$$
  - Computation of traces of  $M^k$  for  $k = jr + i, i = 1, \dots, r - 1, j = 1, \dots, r - 1$ .
    - For  $i$  from 1 to  $r - 1$ 
      - For  $j$  from 1 to  $r - 1$ 

$$N_{jr+i} = \text{Tr}(B_i C_j).$$
  - Computation of the coefficients of  $\det(X \text{Id}_n - M)$ : use Algorithm 8.11 (Newton Sums) taking as  $i$ -th Newton sum  $N_i, i = 0, \dots, n$ .

**Proof of correctness:** Since a square matrix with coefficients in a field  $K$  can be triangulated over  $C$ , the fact that the trace of  $M^i$  is the Newton sums of the eigenvalues is clear in the case of an integral domain. For a general ring with integer division, it is sufficient to specialize the preceding algebraic identity expressed in the ring  $\mathbb{Z}[U_{i,j}, i = 1, \dots, n, j = 1, \dots, n]$  by replacing  $U_{i,j}$  with the entries  $m_{i,j}$  of the matrix. □



**Complexity analysis:** The first step and second step take  $O(rn^3) = O(n^{3.5})$  arithmetic operations. The third step take  $O(n^3)$  arithmetic operations, and the fourth step  $O(n^2)$ .

If the entries of  $M$  are elements of  $\mathbb{Z}$  of bitsize at most  $\tau$ , and the bitsize of  $n$  is  $\nu$ , the bitsizes of the intermediate computations are bounded by  $O((\tau + \nu)n)$  using the complexity analysis of Algorithm 8.14 (Multiplication of Several Matrices) The arithmetic operations performed are multiplications between integers of bitsizes bounded by  $(\tau + \nu)\sqrt{n}$  and integers of bitsizes bounded by  $(\tau + \nu)n$ .

If the entries of  $M$  are elements of  $\mathbb{Z}[Y]$ ,  $Y = Y_1, \dots, Y_k$  of degrees at most  $d$  and of bitsizes at most  $\tau$ , the degrees in  $Y$  and bitsizes of the intermediate computations are  $d n$  and  $(\tau + 2\nu)n$  where  $n$  is the bitsize of  $n d + 1$  using the complexity analysis of Algorithm 8.14 (Multiplication of Several Matrices). The arithmetic operations performed are multiplications between integers of bitsizes bounded by  $(\tau + 2\nu)\sqrt{n}$  and integers of bitsizes bounded by  $(\tau + 2\nu)n$ .  $\square$

*Remark 8.19.* a) In the case of a field of characteristic zero, the rank of  $M$  is easily computed from its characteristic polynomial  $\text{CharPol}(M)$ : it is the degree of the monomial of least degree in  $\text{CharPol}(M)$ .

b) Algorithm 8.17 (Characteristic polynomial) provides the determinant of  $M$  in  $O(n^{3.5})$  arithmetic operations in an arbitrary ring with integer division, substituting 0 to  $X$  in  $\text{CharPol}(M)$ .  $\square$

### 8.2.4 Signature of Quadratic Forms

A general method for computing the signature of quadratic form using the characteristic polynomial is based on the following result.

**Proposition 8.20.** *If  $\Phi$  is a quadratic form with associated symmetric matrix  $M$  of size  $n$ , with entries in a real closed field  $\mathbb{R}$  and*

$$\text{CharPol}(M) = \det(X \text{Id}_n - M) = X^n + a_{n-1}X^{n-1} + \dots + a_0$$

*is the characteristic polynomial of  $M$ , then*

$$\text{Sign}(M) = \text{Var}(1, a_{n-1}, \dots, a_0) - \text{Var}((-1)^n, (-1)^{n-1}a_{n-1}, \dots, a_0),$$

*(see Notation 2.32).*

**Proof:** All the roots of the characteristic polynomial of a symmetric matrix belong to  $\mathbb{R}$  by Theorem 4.42 and we can apply Proposition 2.33 (Descartes' law of signs) and Remark 2.38.  $\square$

**Algorithm 8.18. [Signature Through Descartes]**

- **Structure:** an ordered integral domain  $D$ .
- **Input:** an  $n \times n$  symmetric matrix  $M = [m_{i,j}]$ , with coefficients in  $D$ .
- **Output:** the signature of the quadratic form associated to  $M$ .
- **Complexity:**  $O(n^{3.5})$ .
- **Procedure:** Compute the characteristic polynomial of  $M$

$$\text{CharPol}(M) = \det(X \text{Id}_n - M) = X^n + a_{n-1}X^{n-1} + \dots + a_0$$

using Algorithm 8.17 (Characteristic polynomial) and output

$$\text{Var}(1, a_{n-1}, \dots, a_0) - \text{Var}((-1)^n, (-1)^{n-1}a_{n-1}, \dots, a_0).$$

**Complexity analysis:** The complexity is bounded by  $O(n^{3.5})$ , according to the complexity analysis of Algorithm 8.17 (Characteristic polynomial). Moreover, if the entries of  $A$  are elements of  $\mathbb{Z}$  of bitsize at most  $\tau$ , the arithmetic operations performed are multiplications between integers of bitsizes bounded by  $\tau$  and integers of bitsizes bounded by  $(\tau + 2\nu)n + \nu + 2$  where  $\nu$  is the bitsize of  $n$ .  $\square$

## 8.3 Remainder Sequences and Subresultants

### 8.3.1 Remainder Sequences

We now present some results concerning the computation of the signed remainder sequence that was defined in Chapter 1 (Definition 1.2).

The following algorithm follows immediately from the definition.

**Algorithm 8.19. [Signed Remainder Sequence]**

- **Structure:** a field  $K$ .
- **Input:** two univariate polynomials  $P$  and  $Q$  with coefficients  $K$ .
- **Output:** the signed remainder sequence of  $P$  and  $Q$ .
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:**
  - Initialization:  $i := 1$ ,  $\text{SRemS}_0(P, Q) := P$ ,  $\text{SRemS}_1(P, Q) := Q$ .
  - While  $\text{SRemS}_i(P, Q) \neq 0$ 
    - $\text{SRemS}_{i+1}(P, Q) = -\text{Rem}(\text{SRemS}_{i-1}(P, Q), \text{SRemS}_i(P, Q))$ ,
    - $i := i + 1$ .

**Complexity analysis:** Let  $P$  and  $Q$  have degree  $p$  and  $q$ . The number of steps in the algorithm is at most  $q + 1$ . Denoting by  $d_i = \deg(\text{SRemS}_i(P, Q))$ , the complexity of computing  $\text{SRemS}_{i+1}(P, Q)$  knowing  $\text{SRemS}_{i-1}(P, Q)$  and  $\text{SRemS}_i(P, Q)$  is bounded by  $(d_{i-1} - d_i + 1)(2d_i + 3)$  by Algorithm 8.3. Summing over all  $i$  and bounding  $d_i$  by  $q$ , we get the bound  $(p + q + 1)(2q + 3)$ , which is  $O(pq)$ .  $\square$

An important variant of Signed Euclidean Division is the following Extended Signed Euclidean Division computing the extended signed remainder sequence (Definition 1.10).

*Algorithm 8.20. [Extended Signed Remainder Sequence]*

- **Structure:** a field  $K$ .
- **Input:** two univariate polynomials  $P$  and  $Q$  with coefficients in  $K$ .
- **Output:** the extended signed remainder sequence  $\text{Ex}(P, Q)$ .
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:**
  - Initialization:  $i := 1$ ,

$$\text{SRemS}_0(P, Q) := P,$$

$$\text{SRemS}_1(P, Q) := Q,$$

$$\text{SRemU}_0(P, Q) = \text{SRemV}_1(P, Q) := 1,$$

$$\text{SRemV}_0(P, Q) = \text{SRemU}_1(P, Q) := 0.$$

- While  $\text{SRemS}_i(P, Q) \neq 0$ 
  - Compute

$$A_{i+1} = \text{Quo}(\text{SRemS}_{i-1}(P, Q), \text{SRemS}_i(P, Q)),$$

$$\text{SRemS}_{i+1}(P, Q) = -\text{SRemS}_{i-1}(P, Q) + A_{i+1}\text{SRemS}_i(P, Q),$$

$$\text{SRemU}_{i+1}(P, Q) = -\text{SRemU}_{i-1}(P, Q) + A_{i+1}\text{SRemU}_i(P, Q),$$

$$\text{SRemV}_{i+1}(P, Q) = -\text{SRemV}_{i-1}(P, Q) + A_{i+1}\text{SRemV}_i(P, Q).$$

- $\text{Ex}_i(P, Q) = (\text{SRemS}_i(P, Q), \text{SRemU}_i(P, Q), \text{SRemV}_i(P, Q))$
- $i := i + 1$ .

**Proof of correctness:** Immediate by Proposition 1.9. □

**Complexity analysis:** Suppose that  $P$  and  $Q$  have respective degrees  $p$  and  $q$ . It is immediate to check that the complexity is  $O(pq)$ , as in Algorithm 8.19 (Signed Remainder Sequence). □

If we also take into consideration the growth of the bitsizes of the coefficients in the signed remainder sequence, an exponential behavior of the preceding algorithms is a priori possible. If the coefficients are integers of bitsize  $\tau$ , the bitsizes of the coefficients in the signed remainder sequence of  $P$  and  $Q$  could be exponential in the degrees of the polynomials  $P$  and  $Q$  since the bitsize of the coefficients could be doubled at each computation of a remainder in the euclidean remainder sequence.

The bitsizes of the coefficients in the signed remainder sequence can indeed increase dramatically as we see in the next example.

*Example 8.21.* Consider the following numerical example:

$$\begin{aligned} P := & 9X^{13} - 18X^{11} - 33X^{10} + 102X^8 + 7X^7 - 36X^6 \\ & - 122X^5 + 49X^4 + 93X^3 - 42X^2 - 18X + 9. \end{aligned}$$

The greatest common divisor of  $P$  and  $P'$  is of degree 5. The leading coefficients of the signed remainder sequence of  $P$  and  $P'$  are:

$$\begin{array}{r}
 36 \\
 \hline
 13, \\
 - \frac{10989}{16}, \\
 \hline
 2228672 \\
 - \frac{165649}{900202097355}, \\
 \hline
 4850565316, \\
 - \frac{3841677139249510908}{543561530761725025}, \\
 \hline
 6648854900739944448789496725 \\
 - \frac{676140352527579535315696712}{200117670554781699308164692478544184}, \\
 \hline
 1807309302290980501324553958871415645.
 \end{array}
 \quad \square$$

### 8.3.2 Signed Subresultant Polynomials

Now we define and study the subresultant polynomials. Their coefficients are determinants extracted from the Sylvester matrix, and they are closely related to the remainder sequence. Their coefficients of highest degree are the subresultant coefficients introduced in Chapter 4 and used to study the geometry of semi-algebraic sets in Chapter 5. We are going to use them in this chapter to estimate the bitsizes of the coefficients in the signed remainder sequence. They will be also used for real root counting with a good control on the size of the intermediate computations.

#### 8.3.2.1 Polynomial Determinants

We first study polynomial determinants, which will be useful in the study of subresultants.

Let  $K$  be a field of characteristic 0. Consider the  $K$ -vector space  $\mathcal{F}_n$ , consisting of polynomials whose degrees are less than  $n$ , equipped with the basis

$$\mathcal{B} = X^{n-1}, \dots, X, 1.$$

We associate to a list of polynomials  $\mathcal{P} = P_1, \dots, P_m$ , with  $m \leq n$  a matrix  $\text{Mat}(\mathcal{P})$  whose rows are the coordinates of the  $P_i$ 's in the basis  $\mathcal{B}$ . Note that  $\text{Mat}(\mathcal{B})$  is the identity matrix of size  $n$ .

Let  $0 < m \leq n$ . A mapping  $\Phi$  from  $(\mathcal{F}_n)^m$  to  $\mathcal{F}_{n-m+1}$  is **multilinear** if for  $\lambda \in K, \mu \in K$

$$\Phi(\dots, \lambda A_i + \mu B_i, \dots) = \lambda \Phi(\dots, A_i, \dots) + \mu \Phi(\dots, B_i, \dots).$$

A mapping  $\Phi$  from  $(\mathcal{F}_n)^m$  to  $\mathcal{F}_{n-m+1}$  is **alternating** if

$$\Phi(\dots, A, \dots, A, \dots) = 0.$$

A mapping  $\Phi$  from  $(\mathcal{F}_n)^m$  to  $\mathcal{F}_{n-m+1}$  is **antisymmetric** if

$$\Phi(\dots, A, \dots, B, \dots) = -\Phi(\dots, B, \dots, A, \dots).$$

**Lemma 8.22.** *A mapping from  $(\mathcal{F}_n)^m$  to  $\mathcal{F}_{n-m+1}$  which is multilinear and alternating is antisymmetric.*

**Proof:** Since  $\Phi$  is alternating,

$$\begin{aligned} \Phi(\dots, A+B, \dots, A+B, \dots) &= \Phi(\dots, A, \dots, A, \dots) \\ &= \Phi(\dots, B, \dots, B, \dots) \\ &= 0. \end{aligned}$$

Using multilinearity, we get easily

$$\Phi(\dots, A, \dots, B, \dots) + \Phi(\dots, B, \dots, A, \dots) = 0. \quad \square$$

**Proposition 8.23.** *There exists a unique multilinear alternating mapping  $\Phi$  from  $(\mathcal{F}_n)^m$  to  $\mathcal{F}_{n-m+1}$  satisfying, for every  $n > i_1 > \dots > i_{m-1} > i$*

$$\begin{cases} \Phi(X^{i_1}, \dots, X^{i_{m-1}}, X^i) = X^i & \text{if for every } j < m \text{ } i_j = n - j. \\ \Phi(X^{i_1}, \dots, X^{i_{m-1}}, X^i) = 0 & \text{otherwise.} \end{cases}$$

**Proof:** Decomposing each  $P_i$  in the basis  $\mathcal{B}$  of monomials and using multilinearity and antisymmetry, it is clear that a multilinear and alternating mapping  $\Phi$  from  $\mathcal{F}_n^m$  to  $\mathcal{F}_{n-m+1}$  depends only on the values  $\Phi(X^{i_1}, \dots, X^{i_{m-1}}, X^i)$  for  $n > i_1 > \dots > i_{m-1} > i$ . This proves the uniqueness.

In order to prove existence, let  $m_i, i \leq n$ , be the  $m \times m$  minor of  $\text{Mat}(\mathcal{P})$  based on the columns  $1, \dots, m-1, n-i$ , then

$$\Phi(\mathcal{P}) = \sum_{i \leq n-m} m_i X^i \quad (8.10)$$

satisfies all the properties required. □

The  $(m, n)$ -**polynomial determinant** mapping, denoted  $\text{pdet}_{m,n}$ , is the unique multilinear alternating mapping from  $\mathcal{F}_n^m$  to  $\mathcal{F}_{n-m+1}$  satisfying the properties of Proposition 8.23.

When  $n = m$ , it is clear that  $\text{pdet}_{n,n}(\mathcal{P}) = \det(\text{Mat}(\mathcal{P}))$ , since  $\det$  is known to be the unique multilinear alternating map sending the identity matrix to 1.

On the other hand, when  $m = 1$ ,  $\text{pdet}(P)_{1,n}(X^i) = X^i$  and, by linearity,  $\text{pdet}_{1,n}(P) = P$ .

It follows immediately from the definition that

**Lemma 8.24.** *Let  $\mathcal{P} = P_1, \dots, P_m$ .*

*If  $\mathcal{Q} = Q_1, \dots, Q_m$  is such that  $Q_i = P_i, i \neq j, Q_j = P_j + \sum_{j \neq i} \lambda_j P_j$ , then  $\text{pdet}_{m,n}(\mathcal{Q}) = \text{pdet}_{m,n}(\mathcal{P})$ .*

If  $\mathcal{Q} = P_m, \dots, P_1$ , then  $\text{pdet}_{n,m}(\mathcal{Q}) = \varepsilon_m \text{pdet}_{m,n}(\mathcal{P})$ , where  $\varepsilon_m = (-1)^{m(m-1)/2}$  (see Notation 4.26).

We consider now a sequence  $\mathcal{P}$  of polynomials with coefficients in a ring  $D$ . Equation (8.10) provides a definition of the  $(m, n)$ -polynomial determinant  $\text{pdet}_{m,n}(\mathcal{P})$  of  $\mathcal{P}$ . Note that  $\text{pdet}_{m,n}(\mathcal{P}) \in D[X]$ .

We can express the polynomial determinant as the classical determinant of a matrix whose last column has polynomial entries in the following way:

If  $\mathcal{P} = P_1, \dots, P_m$  we let  $\text{Mat}(\mathcal{P})^*$  be the  $m \times m$  matrix whose first  $m - 1$  columns are the first  $m - 1$  columns of  $\text{Mat}(\mathcal{P})$  and such that the elements of the last column are the polynomials  $P_1, \dots, P_m$ .

With this notation, we have

**Lemma 8.25.**

$$\text{pdet}_{m,n}(\mathcal{P}) = \det(\text{Mat}(\mathcal{P})^*).$$

**Proof:** Using the linearity of  $\det(\text{Mat}(\mathcal{P})^*)$  as a function of its last column, it is clear that  $\det(\text{Mat}(\mathcal{P})^*) = \sum_{i \leq n} m_i X^i$ , using the notation of Proposition 8.23. For  $i > n - m$ ,  $m_i = 0$  since it is the determinant of a matrix with two equal columns. □

*Remark 8.26.* Expanding  $\det(\text{Mat}(\mathcal{P})^*)$  by its last column we observe that  $\text{pdet}_{m,n}(\mathcal{P})$  is a linear combination of the  $P_i$  with coefficients equal (up to sign)  $(m - 1) \times (m - 1)$  to minors extracted on the  $m - 1$  first columns of  $\mathcal{P}$ . It is thus a linear combination with coefficients in  $D$  of the  $P_i$ 's. □

The following immediate consequences of Lemma 8.25 will be useful.

**Lemma 8.27.** *Let  $\mathcal{P} = P_1, \dots, P_\ell, P_{\ell+1}, \dots, P_m$  be such that*

$$\deg(P_i) = n - i, i \leq \ell, \deg(P_i) < n - 1 - \ell, \ell < i \leq m,$$

*with*

$$\begin{aligned} P_i &= p_{i,n-i} X^{n-i} + \dots + p_{i,0}, i \leq \ell, \\ P_i &= p_{i,n-1-\ell} X^{n-1-\ell} + \dots + p_{i,0}, \ell < i \leq m. \end{aligned}$$

*Then*

$$\text{pdet}_{m,n}(\mathcal{P}) = \prod_{i=1}^{\ell} p_{i,n-i} \text{pdet}_{m-\ell,n-\ell}(\mathcal{Q}),$$

*where  $\mathcal{Q} = P_{\ell+1}, \dots, P_m$ .*

**Proof:** Let

$$\begin{aligned} P_i &= p_{i,n-i} X^{n-i} + \dots + p_{i,0}, i \leq \ell, \\ P_i &= p_{i,n-1-\ell} X^{n-1-\ell} + \dots + p_{i,0}, \ell < i \leq m. \end{aligned}$$

The shape of the matrix  $\text{Mat}(\mathcal{P})$  is as follows

$$\begin{bmatrix} p_{1,n-1} & \cdots & \cdots & \cdots & \cdots & \cdots & p_{1,0} \\ 0 & \ddots & & & & & \\ \vdots & \ddots & \ddots & & & & \\ \vdots & & \ddots & p_{\ell,n-\ell} & \cdots & \cdots & p_{\ell,0} \\ \vdots & & & 0 & p_{\ell+1,n-\ell-1} & \cdots & p_{\ell+1,0} \\ \vdots & & & \vdots & \vdots & & \vdots \\ 0 & \cdots & \cdots & 0 & p_{m,n-\ell-1} & \cdots & p_{m,0} \end{bmatrix}$$

Using Lemma 8.25, develop the determinant  $\det(\text{Mat}(\mathcal{P})^*)$  by its first  $\ell$  columns. □

**Lemma 8.28.** *Let  $\mathcal{P} = P_1, \dots, P_m$  be such that for every  $i, 1 \leq i \leq m$ , we have  $\deg(P_i) < n - 1$ . Then*

$$\text{pdet}_{m,n}(\mathcal{P}) = 0.$$

**Proof:** Using Lemma 8.25, develop the determinant  $\det(\text{Mat}(\mathcal{P})^*)$  by its first column which is zero. □

### 8.3.2.2 Definition of Signed Subresultants

For the remainder of this chapter, let  $P$  and  $Q$  be two non-zero polynomials of degrees  $p$  and  $q$ , with  $q < p$ , with coefficients in an integral domain  $D$ . The fraction field of  $D$  is denoted by  $K$ . Let

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + a_{p-2} X^{p-2} + \cdots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \cdots + b_0. \end{aligned}$$

We define the signed subresultants of  $P$  and  $Q$  and some related notions.

**Notation 8.29. [Signed subresultant]** For  $0 \leq j \leq q$ , the  $j$ -th signed subresultant of  $P$  and  $Q$ , denoted  $\text{sResP}_j(P, Q)$ , is the polynomial determinant of the sequence of polynomials

$$X^{q-j-1}P, \dots, P, Q, \dots, X^{p-j-1}Q,$$

with associated matrix the Sylvester-Habicht matrix  $\text{SyHa}_j(P, Q)$  (Notation 4.21). Note that  $\text{SyHa}_j(P, Q)$  has  $p + q - 2j$  rows and  $p + q - j$  columns. Clearly,  $\deg(\text{sResP}_j(P, Q)) \leq j$ . By convention, we extend these definitions for  $q < j \leq p$  by

$$\begin{aligned} \text{sResP}_p(P, Q) &= P, \\ \text{sResP}_{p-1}(P, Q) &= Q, \\ \text{sResP}_j(P, Q) &= 0, \quad q < j < p - 1. \end{aligned}$$

Also by convention  $\text{sResP}_{-1}(P, Q) = 0$ . Note that

$$\text{sResP}_q(P, Q) = \varepsilon_{p-q} b_q^{p-q-1} Q. \quad \square$$

The  $j$ -th signed subresultant coefficient of  $P$  and  $Q$   $\text{sRes}_j(P, Q)$ , (Notation 4.21) is the coefficient of  $X^j$  in  $\text{sResP}_j(P, Q)$ ,  $j < p$ .

If  $\deg(\text{sResP}_j(P, Q)) = j$  (equivalently if  $\text{sRes}_j(P, Q) \neq 0$ ) we say that  $\text{sResP}_j(P, Q)$  is **non-defective**. If  $\deg(\text{sResP}_j(P, Q)) = k < j$  we say that  $\text{sResP}_j(P, Q)$  is **defective** of degree  $k$ .

### 8.3.3 Structure Theorem for Signed Subresultants

We are going to see that the non-zero signed subresultants are proportional to the polynomials in the signed remainder sequence. Moreover, the signed subresultant polynomials present the gap structure, graphically displayed by the following diagram: when  $\text{sResP}_{j-1}$  is defective of degree  $k$ ,  $\text{sResP}_{j-1}$  and  $\text{sResP}_k$  are proportional,  $\text{sResP}_{j-2}, \dots, \text{sResP}_{k+1}$  are zero.

The structure theorem for signed subresultants describes precisely this situation. We write  $s_j$  for  $\text{sRes}_j(P, Q)$  and  $t_j$  for  $\text{lcof}(\text{sResP}_j(P, Q))$ . Note that if  $\deg(\text{sResP}_j(P, Q)) = j$ ,  $t_j = s_j$ . In particular  $t_p = s_p = \text{sign}(a_p)$ .

**Theorem 8.30. [Structure theorem for subresultants]** *Let  $0 \leq j < i \leq p + 1$ . Suppose that  $\text{sResP}_{i-1}(P, Q)$  is non-zero and of degree  $j$ .*

- *If  $\text{sResP}_{j-1}(P, Q)$  is zero, then  $\text{sResP}_{i-1}(P, Q) = \text{gcd}(P, Q)$ , and for  $\ell \leq j - 1$ ,  $\text{sResP}_\ell(P, Q)$  is zero.*
- *If  $\text{sResP}_{j-1}(P, Q) \neq 0$  has degree  $k$  then*

$$\begin{aligned} & s_j t_{i-1} \text{sResP}_{k-1}(P, Q) \\ &= -\text{Rem}(s_k t_{j-1} \text{sResP}_{i-1}(P, Q), \text{sResP}_{j-1}(P, Q)). \end{aligned}$$

*If  $j \leq q$ ,  $k < j - 1$ ,  $\text{sResP}_k(P, Q)$  is proportional to  $\text{sResP}_{j-1}(P, Q)$ . More precisely*

$$\begin{aligned} \text{sResP}_\ell(P, Q) &= 0, \quad j - 1 > \ell > k \\ s_k &= \varepsilon_{j-k} \frac{t_{j-1}^{j-k}}{s_j^{j-k-1}} \\ t_{j-1} \text{sResP}_k(P, Q) &= s_k \text{sResP}_{j-1}(P, Q). \end{aligned}$$

(where  $\varepsilon_i = (-1)^{i(i-1)/2}$ )

Note that Theorem 8.30 implies that  $\text{sResP}_{i-1}$  and  $\text{sResP}_j$  are proportional. The following corollary of Theorem 8.30 will be used later in this chapter.

**Corollary 8.31.** *If  $\text{sResP}_{j-1}(P, Q)$  is of degree  $k$ ,*

$$s_j^2 \text{sResP}_{k-1}(P, Q) = -\text{Rem}(s_k t_{j-1} \text{sResP}_j(P, Q), \text{sResP}_{j-1}(P, Q)).$$



**Proof:** Immediate from Theorem 8.30, using

$$s_j t_{i-1} \text{sResP}_{k-1}(P, Q) = -\text{Rem}(s_k t_{j-1} \text{sResP}_{i-1}(P, Q), \text{sResP}_{j-1}(P, Q))$$

and the proportionality between  $\text{sResP}_{i-1}$  and  $\text{sResP}_j$ .  $\square$

Note that we have seen in Chapter 4 (Proposition 4.24) that  $\deg(\gcd(P, Q))$  is the smallest  $j$  such that  $\text{sRes}_j(P, Q) \neq 0$ . The Structure Theorem 8.30 makes this statement more precise:

**Corollary 8.32.** *The last non-zero signed subresultant of  $P$  and  $Q$  is non-defective and a greatest common divisor of  $P$  and  $Q$ .*

**Proof:** Suppose that  $\text{sResP}_j(P, Q) \neq 0$ , and  $\forall \ell < k \text{sResP}_\ell(P, Q) = 0$ . By Theorem 8.30 there exists  $i$  such that  $\deg(\text{sResP}_{i-1}(P, Q)) = j$ , and  $\text{sResP}_{i-1}(P, Q)$  and  $\text{sResP}_j(P, Q)$  are proportional. So  $\text{sResP}_j(P, Q)$  is non-defective and  $\text{sResP}_{i-1}(P, Q)$  is a greatest common divisor of  $P$  and  $Q$ , again by Theorem 8.30.  $\square$

Moreover, a consequence of the Structure Theorem 8.30 is that signed subresultant polynomials are closely related to the polynomials in the signed remainder sequence.

In the non-defective case, we have:

**Corollary 8.33.** *When all  $\text{sResP}_j(P, Q)$  are non-defective,  $j = p, \dots, 0$ , the signed subresultant polynomials are proportional up to a square to the polynomials in the signed remainder sequence.*

**Proof:** We consider the signed remainder sequence

$$\begin{aligned} \text{SRemS}_0(P, Q) &= P, \\ \text{SRemS}_1(P, Q) &= Q, \\ &\vdots \\ \text{SRemS}_{\ell+1}(P, Q) &= -\text{Rem}(\text{SRemS}_{\ell-1}(P, Q), \text{SRemS}_\ell(P, Q)), \\ &\vdots \\ \text{SRemS}_p(P, Q) &= -\text{Rem}(\text{SRemS}_{p-2}(P, Q), \text{SRemS}_{p-1}(P, Q)), \\ \text{SRemS}_{p+1}(P, Q) &= 0, \end{aligned}$$

and prove by induction on  $\ell$  that  $\text{sResP}_{p-\ell}(P, Q)$  is proportional to  $\text{SRemS}_\ell(P, Q)$ .

The claim is true for  $\ell = 0$  and  $\ell = 1$  by definition of  $\text{sResP}_{p(P, Q)}$  and  $\text{sResP}_{p-1}(P, Q)$ .

Suppose that the claim is true up to  $\ell$ . In the non-defective case, the Structure Theorem 8.30 b) implies

$$\begin{aligned} &s_{p-\ell+1}^2 \text{sResP}_{p-\ell-1}(P, Q) \\ &= -\text{Rem}(s_{p-\ell}^2 \text{sResP}_{p-\ell+1}(P, Q), \text{sResP}_{p-\ell}(P, Q)). \end{aligned}$$

By induction hypothesis,  $\text{sResP}_{p-\ell+1}(P, Q)$  and  $\text{sResP}_{p-\ell}(P, Q)$  are proportional to  $\text{SRemS}_{\ell-1}(P, Q)$  and  $\text{SRemS}_{\ell}(P, Q)$ . Thus, by definition of the signed remainder sequence and by equation (8.12)  $\text{sResP}_{p-\ell-1}(P, Q)$ , is proportional to  $\text{SRemS}_{\ell+1}(P, Q)$ .  $\square$

More generally, the signed subresultants are either proportional to polynomials in the signed remainder sequence or zero.

Let us illustrate this property by an example in the defective case. Let

$$\begin{aligned} P &= X^{11} - X^{10} + 1, \\ P' &= 11X^{10} - 10X^9, \end{aligned}$$

The signed remainder sequence is

$$\begin{aligned} \text{SRemS}_0(P, P') &= X^{11} - X^{10} + 1, \\ \text{SRemS}_1(P, P') &= 11X^{10} - 10X^9 \\ \text{SRemS}_2(P, P') &= \frac{10X^9}{121} - 1, \\ \text{SRemS}_3(P, P') &= -\frac{1331X}{10} + 121, \\ \text{SRemS}_4(P, P') &= \frac{275311670611}{285311670611}. \end{aligned}$$

The non-zero signed subresultant polynomials are the following:

$$\begin{aligned} \text{sResP}_{11}(P, P') &= X^{11} - X^{10} + 1, \\ \text{sResP}_{10}(P, P') &= 11X^{10} - 10X^9, \\ \text{sResP}_9(P, P') &= 10X^9 - 121, \\ \text{sResP}_8(P, P') &= -110X + 100, \\ \text{sResP}_1(P, P') &= 2143588810X - 1948717100, \\ \text{sResP}_0(P, P') &:= -275311670611. \end{aligned}$$

It is easy to check that  $\text{sResP}_8(P, P')$  and  $\text{sResP}_1(P, P')$  are proportional.

**Corollary 8.34.** *If  $\text{SRemS}_{\ell-1}(P, Q)$  and  $\text{SRemS}_{\ell}(P, Q)$  are two successive polynomials in the signed remainder sequence of  $P$  and  $Q$ , of degrees  $d(\ell - 1)$  and  $d(\ell)$ , then  $\text{sResP}_{d(\ell-1)-1}(P, Q)$  and  $\text{sResP}_{d(\ell)}(P, Q)$  are proportional to  $\text{SRemS}_{\ell}(P, Q)$ .*

**Proof:** The proof is by induction on  $\ell$ . Note first that  $P = \text{SRemS}_0$  is proportional to  $\text{sResP}_p$ . The claim is true for  $\ell = 1$  by definition of  $\text{sResP}_p(P, Q)$ ,  $\text{sResP}_{p-1}(P, Q)$ , and  $\text{sResP}_q(P, Q)$ . Suppose that the claim is true up to  $\ell$ . The Structure Theorem 8.30 b) implies (with  $i = d(\ell - 2)$ ,  $j = d(\ell - 1)$ ,  $k = d(\ell)$ ) that  $\text{sResP}_{d(\ell)-1}(P, Q)$  is proportional to

$$\text{Rem}(\text{sResP}_{d(\ell-2)-1}, \text{sResP}_{d(\ell-1)-1})(P, Q).$$

By the induction hypothesis,  $\text{sResP}_{d(\ell-2)-1}(P, Q)$  and  $\text{sResP}_{d(\ell-1)-1}(P, Q)$  are proportional to  $\text{SRemS}_{\ell-1}(p, Q)$  and  $\text{SRemS}_{\ell}(P, Q)$ . It follows that  $\text{sResP}_{d(\ell)-1}(P, Q)$  is proportional to  $\text{SRemS}_{\ell+1}(P, Q)$ . Moreover  $\text{sResP}_{d(\ell)-1}(P, Q)$  and  $\text{sResP}_{d(\ell+1)}(P, Q)$  are proportional by the Structure Theorem 8.30.  $\square$

The proof of the structure theorem relies on the following proposition relating the signed subresultants of  $P$  and  $Q$  and of  $Q$  and  $-R$ , with  $R = \text{Rem}(P, Q)$ .

We recall that  $P$  and  $Q$  are two non-zero polynomials of degrees  $p$  and  $q$ , with  $q < p$ , with coefficients in an integral domain  $D$ , with

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + a_{p-2} X^{p-2} + \dots + a_0, \\ Q &= b_q X^q + b_{q-1} X^{q-1} + \dots + b_0. \end{aligned}$$

The following proposition generalizes Proposition 4.36.

**Proposition 8.35.** *Let  $r$  be the degree of  $R = \text{Rem}(P, Q)$ .*

$$\text{sResP}_j(P, Q) = \varepsilon_{p-q} b_q^{p-r} \text{sResP}_j(Q, -R) \text{ if } j < q-1,$$

where  $\varepsilon_i = (-1)^{i(i-1)/2}$ .

**Proof:** Replacing the polynomials  $X^{q-j-1}P, \dots, P$  by the polynomials  $X^{q-j-1}R, \dots, R$  in  $\text{SyHaPol}_j(P, Q)$  does not modify the polynomial determinant. Indeed,

$$R = P - \sum_{i=0}^{p-q} c_i (X^i Q),$$

where  $C = \sum_{i=0}^{p-q} c_i X^i$  is the quotient of  $P$  in the euclidean division of  $P$  by  $Q$ , and adding to a polynomial of a sequence a multiple of another polynomial of the sequence does not change the polynomial determinant, by Lemma 8.24.

Reversing the order of the polynomials multiplies the polynomial determinant by  $\varepsilon_{p+q-2j}$  using again Lemma 8.24. Replacing  $R$  by  $-R$  multiplies the polynomial determinant by  $(-1)^{q-j}$ , by Lemma 8.24, and  $(-1)^{q-j} \varepsilon_{p+q-2j} = \varepsilon_{p-q}$  (see Notation 4.26). So, defining

$$A_j = \text{pdet}_{p+q-2j, p+q-j}(X^{p-j-1}Q, \dots, Q, -R, \dots, -X^{q-j-1}R),$$

we have

$$\text{sResP}_j(P, Q) = \varepsilon_{p-q} A_j.$$

If  $j \leq r$ ,

$$\begin{aligned} A_j &= b_q^{p-r} \text{pdet}_{q+r-2j, q+r-j}(X^{r-j-1}Q, \dots, Q, -R, \dots, -X^{q-j-1}R) \\ &= b_q^{p-r} \text{sResP}_j(Q, -R), \end{aligned}$$

using Lemma 8.27.

If  $r < j < q - 1$ ,

$$\text{pdet}_{p+q-2j, p+q-j}(X^{p-j-1}Q, \dots, Q, -R, \dots, -X^{q-j-1}R) = 0,$$

using Lemma 8.27 and Lemma 8.28, since  $\deg(-X^{q-j-1}R) < q - 1$ . □

**Proof of Theorem 8.30:** For  $q < j \leq p$ , the only thing to check is that

$$\text{sign}(a_p)^2 \text{sResP}_{q-1}(P, Q) = -\text{Rem}(s_q t_{p-1} \text{sResP}_p(P, Q), \text{sResP}_{p-1}(P, Q)),$$

since  $s_p = \text{sign}(a_p)$  (Notation 4.21) Indeed

$$\begin{aligned} \text{sResP}_{q-1}(P, Q) &= \varepsilon_{p-q} b_q^{p-q+1} R \\ &= -\varepsilon_{p-q+2} b_q^{p-q+1} R \\ &= -\text{Rem}(\varepsilon_{p-q} b_q^{p-q+1} P, Q) \\ &= -\text{Rem}(s_q t_{p-1} P, Q) \end{aligned}$$

since  $s_q = \varepsilon_{p-q} b_q^{p-q}$ ,  $t_{p-1} = b_q$ , and  $\text{sResP}_p(P, Q) = P$ ,  $\text{sResP}_{p-1}(P, Q) = Q$ .

The remainder of the proof is by induction on the length of the remainder sequence of  $P$  and  $Q$ .

Suppose that the theorem is true for  $Q, -R$ . The fact that the theorem holds for  $P, Q$  for  $j \leq r$  is clear by Proposition 8.35, since  $\text{sResP}_j(P, Q)$  and  $\text{sResP}_j(Q, -R)$ ,  $j \leq r$ , are proportional, with the same factor of proportionality  $\varepsilon_{p-q} b_q^{p-r}$ .

For  $r < j \leq q$ , the only thing to check is that

$$s_q t_{p-1} \text{sResP}_{r-1}(P, Q) = -\text{Rem}(s_r t_{q-1} \text{sResP}_{p-1}(P, Q), \text{sResP}_{q-1}(P, Q)),$$

which follows from the induction hypotheses

$$s'_q{}^2 \text{sResP}_{r-1}(Q, -R) = -\text{Rem}(s'_r t'_{q-1} \text{sResP}_q(Q, -R), \text{sResP}_{q-1}(Q, -R)),$$

where we write  $s'_j$  for  $\text{sRes}_j(Q, -R)$  and  $t'_j$  for  $\text{lcof}(\text{sResP}_j(Q, -R))$ , noting that  $t'_q = s'_q = \text{sign}(b - q)$ , since

$$\begin{aligned} &t_{p-1} s_q \text{sResP}_{r-1}(P, Q) \\ &= b_q (\varepsilon_{p-q} b_q^{p-q}) \varepsilon_{p-q} b_q^{p-r} \text{sResP}_{r-1}(Q, -R) \\ &= -b_q^{2p-q-r+1} \text{Rem}(s'_r t'_{q-1} \text{sResP}_q(Q, -R), \text{sResP}_{q-1}(Q, -R)) \\ &= -\text{Rem}(s_r t_{q-1} \text{sResP}_{p-1}(P, Q), \text{sResP}_{q-1}(P, Q)), \end{aligned}$$

by Proposition 8.35, noting that  $t_{q-1} = \varepsilon_{p-q} b_q^{p-q+1} t'_{q-1}$ ,  $s_r = \varepsilon_{p-q} b_q^{p-r} s'_r$  and using that  $\text{sResP}_{q-1}(P, Q)$  is proportional to  $\text{sResP}_{q-1}(Q, -R)$ . □

The following proposition gives a useful precision to Theorem 8.30.

**Proposition 8.36.** *Using the notation of Theorem 8.30 and defining  $C_{k-1}$  as the quotient of  $s_k t_{j-1} \text{sResP}_{i-1}(P, Q)$  by  $\text{sResP}_{j-1}(P, Q)$ , we have  $C_{k-1} \in D[X]$ .*

Before proving it, we need an analogue of Proposition 1.9 for subresultants.

**Notation 8.37. [Subresultant cofactors]** Define  $\text{sResU}_j(P, Q)$  (resp.  $\text{sResV}_j(P, Q)$ ) as  $\det(M_i)$  (resp.  $\det(N_i)$ ), where  $M_i$  (resp.  $N_i$ ) is the square matrix obtained by taking the first  $p+q-2j-1$  columns of  $\text{SyHa}_j(P, Q)$  and with last column equal to  $(X^{q-1-j}, \dots, X, 1, 0, \dots, 0)^t$  (resp.  $(0, \dots, 0, 1, X, \dots, X^{p-1-j})^t$ ).

Note that if  $P, Q \in \mathbb{D}[X]$ , then  $\text{sResU}_j(P, Q), \text{sResV}_j(P, Q) \in \mathbb{D}[X]$ .  $\square$

**Proposition 8.38.** *Let  $j \leq q$ . Then,*

a)  $\deg(\text{sResU}_{j-1}(P, Q)) \leq q-j, \deg(\text{sResV}_{j-1}(P, Q)) \leq p-j,$

$$\text{sResP}_j(P, Q) = \text{sResU}_j(P, Q)P + \text{sResV}_j(P, Q)Q.$$

b) *If  $\text{sResP}_j(P, Q)$  is not 0 and if  $U$  and  $V$  are such that*

$$UP + VQ = \text{sResP}_j(P, Q),$$

*$\deg(U) \leq q-j-1$ , and  $\deg(V) \leq p-j-1$ , then  $U = \text{sResU}_j(P, Q)$  and  $V = \text{sResV}_j(P, Q)$ .*

c) *If  $\text{sResP}_j(P, Q)$  is non-defective, then*

$$\deg(\text{sResU}_{j-1}(P, Q)) = q-j, \deg(\text{sResV}_{j-1}(P, Q)) = p-j,$$

*and  $\text{lcof}(\text{sResV}_{j-1}(P, Q)) = a_p \text{sRes}_j(P, Q)$ .*

**Proof:** a) The conditions

$$\deg(\text{sResU}_{j-1}(P, Q)) \leq q-j, \deg(\text{sResV}_{j-1}(P, Q)) \leq p-j$$

follow from the definitions of  $\text{sResU}_{j-1}(P, Q)$  and  $\text{sResV}_{j-1}(P, Q)$ . By Lemma 8.25,  $\text{sResP}_j(P, Q) = \det(\text{SyHa}_j(P, Q)^*)$ , where  $\text{SyHa}_j(P, Q)^*$  is the square matrix obtained by taking the first  $p+q-2j-1$  columns of  $\text{SyHa}_j(P, Q)$  and with last column equal to

$$(X^{q-1-j}P, \dots, XP, P, Q, \dots, X^{p-j-1}Q)^t.$$

Expanding the determinant by its last column, we obtain the claimed identity.

b) Suppose  $\deg(U) \leq q-j-1, \deg(V) \leq p-j-1$ , and

$$\text{sResP}_j(P, Q) = UP + VQ$$

so that

$$(\text{sResU}_j(P, Q) - U)P + (\text{sResV}_j(P, Q) - V)Q = 0.$$

If  $\text{sResU}_j(P, Q) - U$  is not 0, then  $\text{sResV}_j(P, Q) - V$  cannot be 0, and it follows from Proposition 1.5 that  $\deg(\text{gcd}(P, Q)) > j$ . But this is impossible since  $\text{sResP}_j(P, Q)$  is a non-zero polynomial of degree  $\leq j$  belonging to the ideal generated by  $P$  and  $Q$ .

- c) Since  $\text{sResP}_j(P, Q)$  is non-defective, it follows that  $\text{sRes}_j(P, Q) \neq 0$ . By considering the determinant of the matrix  $\text{SyHa}_{j-1}(P, Q)^*$ , it is clear that the coefficient of  $X^{p-j}$  in  $\text{sResV}_{j-1}(P, Q)$  is  $a_p \text{sRes}_j(P, Q)$ . Moreover,

$$\deg(\text{sResV}_{j-1}) = p - j, \deg(\text{sResU}_{j-1}(P, Q)) = q - j.$$

□

We omit  $P$  and  $Q$  in the notation in the next paragraphs. For  $\text{sResP}_{i-1}$  non-zero of degree  $j$ , we define

$$B_{j,i} = \begin{bmatrix} \text{sResU}_{i-1} & \text{sResV}_{i-1} \\ \text{sResU}_{j-1} & \text{sResV}_{j-1} \end{bmatrix},$$

where  $\text{sResU}_{i-1}, \text{sResV}_{i-1}, \text{sResU}_{j-1}, \text{sResV}_{j-1} \in \mathbb{D}[X]$  are the polynomials of the  $(i-1)$ -th and  $(j-1)$ -th relations of Proposition 8.38, whence

$$\begin{bmatrix} \text{sResP}_{i-1} \\ \text{sResP}_{j-1} \end{bmatrix} = B_{j,i} \cdot \begin{bmatrix} P \\ Q \end{bmatrix}. \quad (8.11)$$

**Lemma 8.39.** *If  $\text{sResP}_{i-1}$  is non-zero of degree  $j$ , then*

$$\det(B_{j,i}) = s_j t_{i-1}.$$

**Proof:** Eliminating  $Q$  from the system (8.11), we have

$$\begin{aligned} & (\text{sResU}_{i-1} \text{sResV}_{j-1} - \text{sResU}_{j-1} \text{sResV}_{i-1}) P \\ &= \text{sResV}_{j-1} \text{sResP}_{i-1} - \text{sResV}_{i-1} \text{sResP}_{j-1}. \end{aligned}$$

Since  $\deg(SR_{i-1}) = j$ ,  $\deg(SR_j) = j$  by the Structure Theorem 8.30, and  $\deg(\text{sResV}_{j-1}) = p - j$ . Using  $\deg(SR_{j-1}) \leq j - 1$  and  $\deg(\text{sResV}_{i-1}) \leq p - i < p - j$ , we see that the right hand side of equation (8.14) has degree  $p$ . The leading coefficient of  $\text{sResV}_{j-1}$  is  $a_p s_j$  by Proposition 8.38. Hence

$$\text{sResU}_{i-1} \text{sResV}_{j-1} - \text{sResU}_{j-1} \text{sResV}_{i-1} = s_j t_{i-1} \neq 0. \quad \square$$

**Corollary 8.40.** *If  $\text{sResP}_{i-1}$  is non-zero of degree  $j$ , then*

$$B_{j,i}^{-1} = \frac{1}{s_j t_{i-1}} \begin{bmatrix} \text{sResV}_{j-1} & -\text{sResV}_{i-1} \\ -\text{sResU}_{j-1} & \text{sResU}_{i-1} \end{bmatrix}, \text{ and } s_j t_{i-1} B_{j,i}^{-1} \in \mathbb{D}[X].$$

Now we study the transition between two consecutive couples of signed subresultant polynomials  $\text{sResP}_{i-1}, \text{sResP}_{j-1}$  and  $\text{sResP}_{j-1}, \text{sResP}_{k-1}$ , where  $\text{sResP}_{i-1}$  is of degree  $j$ ,  $\text{sResP}_{j-1}$  is of degree  $k$ , and  $0 \leq k < j \leq p$ .

The **signed subresultant transition matrix** is

$$T_j = \begin{bmatrix} 0 & 1 \\ -\frac{s_k t_{j-1}}{s_j t_{i-1}} & \frac{C_{k-1}}{s_j t_{i-1}} \end{bmatrix} \in \mathbb{K}[X]^{2 \times 2},$$

so that

$$\text{sResP}_{k-1} = -\frac{s_k t_{j-1}}{s_j t_{i-1}} \text{sResP}_{i-1} + \frac{C_{k-1}}{s_j t_{i-1}} \text{sResP}_{j-1} \tag{8.12}$$

and

$$\begin{bmatrix} \text{sResP}_{j-1} \\ \text{sResP}_{k-1} \end{bmatrix} = T_j \begin{bmatrix} \text{sResP}_{i-1} \\ \text{sResP}_{j-1} \end{bmatrix} \tag{8.13}$$

by the Structure Theorem 8.30.

**Lemma 8.41.** *If  $\text{sResP}_{i-1}$  is non-zero of degree  $j$  and  $\text{sResP}_{j-1}$  is non-zero of degree  $k$ , then*

$$B_{k,j} = T_j B_{j,i}.$$

**Proof:** Let

$$T_j B_{j,i} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}.$$

A simple degree calculation shows that  $\deg(A) \leq q - j$ ,  $\deg(B) \leq p - j$ , and  $\deg(C) = q - k$ , and  $\deg(D) = p - k$ . From equations (8.13) and (8.11) we see that

$$\begin{aligned} \text{sResP}_{j-1} &= AP + BQ \\ \text{sResP}_{k-1} &= CP + DQ. \end{aligned}$$

The conclusion follows from the uniqueness asserted in Proposition 8.38 b).  $\square$

**Proof of Proposition 8.36:** From Lemma 8.41, we see that  $T_j = B_{k,j} B_{j,i}^{-1}$ , which together with the definition of  $B_{k,j}$  and Corollary 8.40 shows that

$$\frac{C_{k-1}}{s_j t_{i-1}} = \frac{1}{s_j t_{i-1}} (-\text{sResU}_{k-1} \text{sResV}_{i-1} + \text{sResV}_{k-1} \text{sResU}_{i-1}),$$

whence  $C_{k-1} = \text{sResU}_{k-1} \text{sResV}_{i-1} - \text{sResV}_{k-1} \text{sResU}_{i-1} \in D[X]$ .  $\square$

**Proposition 8.42.** *Let  $j \leq q$ ,  $\deg(\text{sResP}_j) = j$   $\deg(\text{sResP}_{j-1}) = k \leq j - 1$ ,*

$$\begin{aligned} S_{j-1} &= \text{sResP}_{j-1}, \\ S_{j-1-\delta} &= \frac{(-1)^\delta t_{j-1} S_{j-\delta}}{s_j}, \text{ for } \delta = 1, \dots, j - k - 1. \end{aligned}$$

*Then all of these polynomials are in  $D[X]$  and  $\text{sResP}_k = S_k$ .*

**Proof:** Add the  $j - k - 1 - \delta$  polynomials  $X^{k+\delta+1}, \dots, X^j$  to  $\text{SyHaPol}_{j-1}$  to obtain  $M_{j-1-\delta}$ . It is easy to see that the polynomial determinant of  $M_{j-1-\delta}$  is  $S_{j-1-\delta}$ .  $\square$

### 8.3.4 Size of Remainders and Subresultants

Observe, comparing the following example with Example 8.21, that the bit-sizes of coefficients in the signed subresultant sequence can be much smaller than in the signed remainder sequence.

*Example 8.43.* We consider, as in Example 8.21,

$$P := 9X^{13} - 18X^{11} - 33X^{10} + 102X^8 + 7X^7 - 36X^6 \\ - 122X^5 + 49X^4 + 93X^3 - 42X^2 - 18X + 9.$$

The subresultant coefficients of  $P$  and  $P'$  for  $j$  from 11 to 5 are:

$$\begin{aligned} & 37908 \\ & - 72098829 \\ & - 666229317948 \\ & - 1663522740400320 \\ & - 2181968897553243072 \\ & - 151645911413926622112 \\ & - 165117711302736225120, \end{aligned}$$

the remaining subresultants being 0. □

The difference in bitsizes of coefficients between signed remainder and signed subresultant sequences observed in Example 8.21 and Example 8.43 is a general fact.

First, let us see that the size of subresultants is well controlled. Indeed, using Proposition 8.10 we obtain the following:

**Proposition 8.44. [Size of signed subresultants]** *If  $P$  and  $Q$  have degrees  $p$  and  $q$  and have coefficients in  $\mathbb{Z}$  which have bitsizes at most  $\tau$ , then the bitsizes of the coefficients of  $\text{sResP}_j(P, Q)$  and of  $\text{sResU}_j$  and  $\text{sResV}_j$  are at most  $(\tau + \nu_j)(p + q - 2j)$ , where  $\nu_j$  is the bitsize of  $p + q - 2j$ .*

We also have, using Proposition 8.11,

**Proposition 8.45. [Degree of signed subresultants]** *If  $P$  and  $Q$  have degrees  $p$  and  $q$  and have coefficients in  $\mathbb{R}[Y_1, \dots, Y_k]$  which have degrees  $d$  in  $Y_1, \dots, Y_k$  then the degree of  $\text{sResP}_j(P, Q)$  in  $Y_1, \dots, Y_k$  is at most  $d(p + q - 2j)$ .*

We finally have, using Proposition 8.12,

**Proposition 8.46.** *If  $P$  and  $Q$  have degrees  $p$  and  $q$  and have coefficients in  $\mathbb{Z}[Y_1, \dots, Y_k]$  which have degrees  $d$  in  $Y_1, \dots, Y_k$  of bitsizes  $\tau$ , then the degree of  $\text{sResP}_j(P, Q)$  in  $Y_1, \dots, Y_k$  is at most  $d(p + q - 2j)$ , and the bitsizes of the coefficients of  $\text{sResP}_j(P, Q)$  are at most  $(\tau + \nu)(p + q - 2j) + k\mu$  where  $\nu$  is the bitsize of  $p + q$  and  $\mu$  is the bitsize of  $(p + q)d + 1$ .*

The relationship between the signed subresultants and remainders provides a bound for the bitsizes of the coefficients of the polynomials appearing in the signed remainder sequence.



**Theorem 8.47. [Size of signed remainders]** *If  $P \in \mathbb{Z}[X]$  and  $Q \in \mathbb{Z}[X]$  have degrees  $p$  and  $q < p$  and have coefficients of bitsizes at most  $\tau$ , then the numerators and denominators of the coefficients of the polynomials in the signed remainder sequence of  $P, Q$  have bitsizes at most  $(p+q)(q+1)(\tau+\nu)+\tau$ , where  $\nu$  is the bitsize of  $p+q$ .*

**Proof:** Denote by

$$P = S_0, Q = S_1, S_2, \dots, S_k$$

the polynomials in the signed remainder sequence of  $P$  and  $Q$ . Let  $d_j = \deg(S_j)$ . According to Theorem 8.30  $S_\ell$  is proportional to  $\text{sResP}_{d_{\ell-1}}$ , which defines  $\beta_\ell \in \mathbb{Q}$  such that

$$S_\ell = \beta_\ell \text{sResP}_{d_{\ell-1}}.$$

Consider successive signed remainders of respective degrees  $i = d_{\ell-3}, j = d_{\ell-2}$ , and  $k = d_{\ell-1}$ . According to Theorem 8.30,

$$s_j t_{i-1} \text{sResP}_{k-1} = -\text{Rem}(s_k t_{j-1} \text{sResP}_{i-1}, \text{sResP}_{j-1}),$$

which implies that

$$\beta_\ell = \frac{s_j t_{i-1}}{s_k t_{j-1}} \beta_{\ell-2}$$

since

$$S_\ell = -\text{Rem}(S_{\ell-2}, S_{\ell-1}).$$

Denoting by  $D_\ell$  and  $N_\ell$  the bitsizes of the numerator and denominator of  $\beta_\ell$ , and using Proposition 8.44, we get the estimates

$$\begin{aligned} N_\ell &\leq 2(p+q)(\tau+\nu) + N_{\ell-2}, \\ D_\ell &\leq 2(p+q)(\tau+\nu) + D_{\ell-2}. \end{aligned}$$

Since  $N_0, D_0, N_1$ , and  $D_1$  are bounded by  $\tau$  and  $\ell \leq q+1$ , the claim follows.  $\square$

This quadratic behavior of the bitsizes of the coefficients of the signed remainder sequence is often observed in practice (see Example 8.21).

### 8.3.5 Specialization Properties of Subresultants

Since the signed subresultant is defined as a polynomial determinant which is a multilinear form with respect to its rows, and given the convention for  $\text{sResP}_p$  (see Notation 8.29), we immediately have the following:

Let  $f: D \rightarrow D'$  be a ring homomorphism, and let  $f$  also denote the induced homomorphism from  $f: D[X] \rightarrow D'[X]$ .

**Proposition 8.48.** *Suppose that  $\deg(f(P)) = \deg(P)$ ,  $\deg(f(Q)) = \deg(Q)$ . Then for all  $j \leq p$ ,*

$$\text{sResP}_j(f(P), f(Q)) = f(\text{sResP}_j(P, Q)).$$

Applying this to the ring homomorphism from  $\mathbb{Z}[Y][X]$  to  $\mathbb{R}[X]$  obtained by assigning values  $(y_1, \dots, y_\ell) \in \mathbb{R}^\ell$  to the variables  $(Y_1, \dots, Y_\ell)$ , we see that the signed subresultants after specialization are obtained by specializing the coefficients of the signed subresultants.

*Example 8.49.* Consider, for example, the general polynomial of degree 4:

$$P = X^4 + a X^2 + b X + c.$$

The signed subresultant sequence of  $P$  and  $P'$  is formed by the polynomials (belonging to  $\mathbb{Z}[a, b, c][X]$ )

$$\begin{aligned} \text{sResP}_4(P, P') &= X^4 + a X^2 + b X + c \\ \text{sResP}_3(P, P') &= 4 X^3 + 2 a X + b \\ \text{sResP}_2(P, P') &= -4 (2 a X^2 + 3 b X + 4 c) \\ \text{sResP}_1(P, P') &= 4 ((8 a c - 9 b^2 - 2 a^3) X - a^2 b - 12 b c) \\ \text{sResP}_0(P, P') &= 256 c^3 - 128 a^2 c^2 + 144 a b^2 c + 16 a^4 c - 27 b^4 - 4 a^3 b^2, \end{aligned}$$

which agree, up to squares in  $\mathbb{Q}(a, b, c)$ , with the signed remainder sequence for  $P$  and  $P'$  when there is a polynomial of each degree in the signed remainder sequence (see example 1.15). If  $a = 0$ , the subresultant sequence of the polynomial  $P = X^4 + b X + c$  and  $P'$  is

$$\begin{aligned} \text{sResP}_4(P, P') &= X^4 + b X + c \\ \text{sResP}_3(P, P') &= 4 X^3 + b \\ \text{sResP}_2(P, P') &= -4 (3 b X + 4 c) \\ \text{sResP}_1(P, P') &= -12 b (3 b X + 4 c) \\ \text{sResP}_0(P, P') &= -27 b^4 + 256 c^3, \end{aligned}$$

which is the specialization of the signed subresultant sequence of  $P$  with  $a = 0$ . Comparing this with Example 1.15, we observe that the polynomials in the signed subresultant sequence are multiples of the polynomials in the signed remainder sequence obtained when  $a = 0$ . We also observe the proportionality of  $\text{sResP}_2$  and  $\text{sResP}_1$ , which is a consequence of the Structure Theorem 8.30. □

Note that if  $f: D \rightarrow D'$  is a ring homomorphism such that  $\deg(f(P)) = \deg(P)$ ,  $\deg(f(Q)) < \deg(Q)$ , then for all  $j \leq \deg(f(Q))$

$$f(\text{sResP}_j(P, Q)) = \text{lcof}(f(P))^{\deg(Q) - \deg(f(Q))} \text{sResP}_j(f(P), f(Q)),$$

using Lemma 8.27.

### 8.3.6 Subresultant Computation

We now describe an algorithm for computing the subresultant sequence, based upon the preceding results.

Let  $P$  and  $Q$  be polynomials in  $D[X]$  with  $\deg(P) = p$ ,  $\deg(Q) = q < p$ . The **signed subresultant sequence** is the sequence

$$\text{sResP}(P, Q) = \text{sResP}_p(P, Q), \dots, \text{sResP}_0(P, Q).$$

*Algorithm 8.21. [Signed Subresultant]*

- **Structure:** an ordered integral domain  $D$ .
- **Input:** two univariate polynomials

$$P = a_p X^p + \dots + a_0$$

$$Q = b_q X^q + \dots + b_0$$

with coefficients  $D$  of respective degrees  $p$  and  $q$ ,  $p > q$ .

- **Output:** the sequence of signed subresultant polynomials and signed subresultant coefficients.
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:**
  - Initialize:

$$\text{sResP}_p := P,$$

$$s_p = t_p := \text{sign}(a_p),$$

$$\text{sResP}_{p-1} := Q,$$

$$t_{p-1} := b_q,$$

$$\text{sResP}_q := \varepsilon_{p-q} b_q^{p-q-1} Q,$$

$$s_q := \varepsilon_{p-q} b_q^{p-q},$$

$$\text{sResP}_\ell = s_\ell := 0 \quad \text{for } \ell \text{ from } q+1 \text{ to } p-2$$

$$i := p+1, j := p.$$

- While  $\text{sResP}_{j-1} \neq 0$ ,
  - $k := \deg(\text{sResP}_{j-1})$ ,
    - If  $k = j-1$ ,
      - $s_{j-1} := t_{j-1}$ .
      - $\text{sResP}_{k-1} := -\text{Rem}(s_{j-1}^2 \text{sResP}_{i-1}, \text{sResP}_{j-1}) / (s_j t_{i-1})$ .
    - If  $k < j-1$ ,
      - $s_{j-1} := 0$ .
      - Compute  $s_k$  and  $\text{sResP}_k$ : for  $\delta$  from 1 to  $j-k-1$ :
 
$$t_{j-\delta-1} := (-1)^\delta (t_{j-1} t_{j-\delta}) / s_j,$$

$$s_k := t_k$$

$$\text{sResP}_k := s_k \text{sResP}_{j-1} / t_{j-1}.$$
    - Compute  $s_\ell$  and  $\text{sResP}_\ell$  for  $\ell$  from  $j-2$  to  $k+1$ :
 
$$\text{sResP}_\ell = s_\ell := 0.$$
  - Compute  $\text{sResP}_{k-1}$ :
 
$$\text{sResP}_{k-1} := -\text{Rem}(t_{j-1} s_k \text{sResP}_{i-1}, \text{sResP}_{j-1}) / (s_j t_{i-1}).$$

- $t_{k-1} := \text{lcof}(\text{sResP}_{k-1})$ .
- $i := j, j := k$ .
- For  $\ell = 0$  to  $j - 2$

$$\text{sResP}_\ell = s_\ell := 0.$$

- Output  $\text{sResP} := \text{sResP}_p, \dots, \text{sResP}_0, \text{sRes} := s_p, \dots, s_0$ .

**Proof of correctness:** The correctness of the algorithm follows from Theorem 8.30.  $\square$

**Complexity analysis:** All the intermediate results in the computation belong to  $D[X]$  by the definition of the signed subresultants as polynomial determinants (Notation 8.29) and Proposition 8.42.

The computation of  $\text{sResP}_{k-1}$  takes  $j + 2$  multiplications to compute  $s_k t_{j-1} \text{sResP}_{i-1}$ ,  $(j - k + 1)(2k + 3)$  arithmetic operations to perform the euclidean division of  $s_k t_{j-1} \text{sResP}_{i-1}$  by  $\text{sResP}_{j-1}$ , one multiplication and  $k$  divisions to obtain the result. The computation of  $s_k$  takes  $j - k - 1$  multiplications and  $j - k - 1$  exact divisions. The computation of  $\text{sResP}_k$  takes  $k + 1$  multiplications and  $k + 1$  exact divisions. So computing  $\text{sResP}_{k-1}$  and  $\text{sResP}_k$  takes  $O((j - k)k)$  arithmetic operations.

Finally the complexity of computing the signed subresultant sequence is  $O(pq)$ , similarly to the computation of the signed remainder sequence when  $q < p$  (Algorithm 8.19).

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , with coefficients of bitsizes bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $(\tau + \nu)(p + q)$  where  $\nu$  is the bitsize of  $p + q$  according to Proposition 8.44.  $\square$

*Remark 8.50.* Note that initializing  $s_p = t_p = 1$  Algorithm 8.21 (Signed Subresultant) is also valid in a domain, and computes correct values of

$$\text{sResP}_{p-1}, \dots, \text{sResP}_0, \text{sRes}_{p-1}, \dots, \text{sRes}_0. \quad \square$$

Note that Algorithm 8.21 (Signed Subresultant) provides an algorithm for computing the resultant of two polynomial of degree  $p$  and  $q$ ,  $q < p$ , with complexity  $O(pq)$ , since  $\text{sRes}_0(P, Q)$  is up to a sign equal to the resultant of  $P$  and  $Q$ , while a naive computation of the resultant as a determinant would have complexity  $O(p^3)$ . This improvement is due to the special structure of the Sylvester-Habicht matrix, which is taken into account in the subresultant algorithm. Algorithm 8.21 (Signed Subresultant) can be used to compute the resultant with complexity  $O(pq)$  in the special case  $p = q$  as well.

**Exercise 8.2.** Describe an algorithm computing the resultant of  $P$  and  $Q$  with complexity  $O(p^2)$  when  $\deg(P) = \deg(Q) = p$ . Hint: consider  $Q_1 = a_p Q - b_p P$  and prove that  $a_p^{p-1} \text{Res}(P, Q) = \text{Res}(P, Q_1)$ .

The signed subresultant coefficients are also computed in time  $O(pq)$  using Algorithm 8.21 (Signed Subresultant), while computing them from their definition as determinants using Algorithm 8.16 (Dogdson-Jordan-Bareiss) would cost  $O(p^4)$ , since there are  $O(p)$  determinants of matrices of size  $O(p)$  to compute.

*Algorithm 8.22. [Extended Signed Subresultant]*

- **Structure:** an ordered integral domain  $D$ .
- **Input:** two univariate polynomials

$$P = a_p X^p + \dots + a_0$$

$$Q = b_q X^q + \dots + b_0$$

with coefficients  $D$  of respective degrees  $p$  and  $q$ ,  $p > q$ .

- **Output:** the sequence of signed subresultant polynomials and the corresponding  $\text{sResU}$  and  $\text{sResV}$ .
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:**
  - Initialize:

$$\text{sResP}_p := P,$$

$$s_p = t_p := \text{sign}(a_p),$$

$$\text{sResP}_{p-1} := Q,$$

$$t_{p-1} := b_q,$$

$$\text{sResU}_p = \text{sResV}_{p-1} := 1,$$

$$\text{sResV}_p = \text{sResU}_{p-1} := 0,$$

$$\text{sResP}_q := \varepsilon_{p-q} b_q^{p-q-1} Q,$$

$$s_q := \varepsilon_{p-q} b_q^{p-q},$$

$$\text{sResU}_q := 0,$$

$$\text{sResV}_q := \varepsilon_{p-q} b_q^{p-q},$$

$$\text{sResP}_\ell = s_\ell = \text{sResU}_\ell = \text{sResV}_\ell := 0 \quad \text{for } \ell \text{ from } q+1 \text{ to } p-2$$

$$i := p+1, \quad j := p.$$

- While  $\text{sResP}_{j-1} \neq 0$ 
  - $k := \deg(\text{sResP}_{j-1})$
  - If  $k = j-1$ ,

$$s_{j-1} := t_{j-1},$$

$$C_{k-1} := \text{Quo}(s_{j-1}^2 \text{sResP}_{i-1}, \text{sResP}_{j-1}),$$

$$\text{sResP}_{k-1} := (-s_{j-1}^2 \text{sResP}_{i-1} + C_{k-1} \text{sResP}_{j-1}) / (s_j t_{i-1}),$$

$$\text{sResU}_{k-1} := (-s_{j-1}^2 \text{sResU}_{i-1} + C_{k-1} \text{sResU}_{j-1}) / (s_j t_{i-1}),$$

$$\text{sResV}_{k-1} := (-s_{j-1}^2 \text{sResV}_{i-1} + C_{k-1} \text{sResV}_{j-1}) / (s_j t_{i-1}).$$

- If  $k < j - 1$ ,

$$s_{j-1} := 0.$$

Compute  $\text{sResP}_k, \text{sResU}_k, \text{sResV}_k$ : for  $\delta$  from 1 to  $j - k - 1$

$$t_{j-\delta-1} := (-1)^\delta (t_{j-1} t_{j-\delta}) / s_j,$$

$$s_k := t_k,$$

$$\text{sResP}_k := s_k \text{sResP}_{j-1} / t_{j-1},$$

$$\text{sResU}_k := s_k \text{sResU}_{j-1} / t_{j-1},$$

$$\text{sResV}_k := s_k \text{sResV}_{j-1} / t_{j-1}.$$

Compute  $s_\ell, \text{sResP}_\ell, \text{sResU}_\ell, \text{sResV}_\ell$ : for  $\ell$  from  $j - 2$  to  $k + 1$ :

$$\text{sResP}_\ell = s_\ell = \text{sResU}_\ell = \text{sResV}_\ell := 0.$$

Compute  $\text{sResP}_{k-1}, \text{sResU}_{k-1}, \text{sResV}_{k-1}$ :

$$C_{k-1} := \text{Quo}(s_k t_{j-1} \text{sResP}_{i-1}, \text{sResP}_{j-1}),$$

$$\text{sResP}_{k-1} := (-s_k t_{j-1} \text{sResP}_{i-1} + C_{k-1} \text{sResP}_{j-1}) / (s_j t_{i-1}),$$

$$\text{sResU}_{k-1} := (-s_k t_{j-1} \text{sResU}_{i-1} + C_{k-1} \text{sResU}_{j-1}) / (s_j t_{i-1}),$$

$$\text{sResV}_{k-1} := (-s_k t_{j-1} \text{sResV}_{i-1} + C_{k-1} \text{sResV}_{j-1}) / (s_j s_{i-1}).$$

- $t_{k-1} := \text{lcof}(\text{sResP}_{k-1}).$

- $i := j, j := k.$

- For  $\ell = j - 2$  to 0:

$$\text{sResP}_\ell = s_\ell = \text{sResU}_\ell = \text{sResV}_\ell := 0.$$

- Output  $\text{sResP} := \text{sResP}_p, \dots, \text{sResP}_0, \text{sResU} := \text{sResU}_p, \dots, \text{sResU}_0,$   
 $\text{sResV} := \text{sResV}_p, \dots, \text{sResV}_0.$

**Proof of correctness:** The correctness of the algorithm follows from Theorem 8.30 and Proposition 8.38 b) since it is immediate to verify that, with  $\text{sResU}$  and  $\text{sResV}$  computed in the algorithm above,

$$\text{sResP}_{i-1} = \text{sResU}_{i-1} P + \text{sResV}_{i-1} Q,$$

$$\text{sResP}_{j-1} = \text{sResU}_{j-1} P + \text{sResV}_{j-1} Q.$$

This implies that  $\text{sResP}_{k-1} = \text{sResU}_{i-1} P + \text{sResV}_{i-1} Q.$  □

**Complexity analysis:** The complexity is clearly  $O(pq)$  as in Algorithm 8.21 (Signed Subresultant).

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , with coefficients of bitsizes bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $(\tau + \nu)(p + q)$ , where  $\nu$  is the bitsize of  $p + q$  according to Proposition 8.44. □

*Remark 8.51.* Algorithm 8.21 (Signed Subresultant) and Algorithm 8.22 (Extended Signed Subresultant) use exact divisions and are valid only in an integral domain, and not in a general ring. In a ring with division by integers, the algorithm computing determinants indicated in Remark 8.19 can always be used for computing the signed subresultant coefficients. The complexity obtained is  $(pq)^{O(1)}$  arithmetic operations in the ring  $D$  of coefficients of  $P$  and  $Q$ , which is sufficient for the complexity estimates obtained in later chapters.  $\square$

## 8.4 Bibliographical Notes

Bounds on determinants are due to Hadamard [80]. A variant of Dogdson-Jordan-Bareiss's algorithm appears in [54] (see also [9]). Note that Dogdson is better known as Lewis Carrol.

The idea of using the traces of the powers of the matrix for computing its characteristic polynomial is a classical method due to Leverrier [106]. The improvement we present is due to Preparata and Sarwate [132]. It is an instance of the "baby step -giant step" method.

Subresultant polynomials and their connection with remainders were already known to Euler [56] and have been studied by Habicht [79]. Subresultants appear in computer algebra with Collins [44], and they have been studied extensively since then.

There are much more sophisticated algorithms than the ones presented in this book, and with much better complexity, for polynomial and matrix multiplication (see von zur Gathen and Gerhard's *Modern Computer Algebra* [64]).

---

## Cauchy Index and Applications

In Section 9.1, several real root and Cauchy index counting methods are described. Section 9.2 deals with the closely related topic of Hankel matrices and quadratic forms. In Section 9.3 an important application of Cauchy index to counting complex roots with positive real part is described. The only ingredient used in later chapters of the book coming from Chapter 9 is the computation of the Tarski-query.

### 9.1 Cauchy Index

#### 9.1.1 Computing the Cauchy Index

A first algorithm for computing the Cauchy index follows from Algorithm 8.19 (Signed Remainder Sequence), using Theorem 2.58 (Sturm).

*Algorithm 9.1. [Sturm Cauchy Index]*

- **Structure:** an ordered field  $K$ .
- **Input:**  $P \in K[X] \setminus \{0\}$ ,  $Q \in K[X]$ .
- **Output:** the Cauchy index  $\text{Ind}(Q/P)$ .
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:** Compute the signed remainder sequence of  $P$  and  $Q$ , using Algorithm 8.19, then compute the difference in sign variations at  $-\infty$  and  $+\infty$  from the degrees and signs of leading coefficients of the polynomials in this sequence.

**Proof of correctness:** The correctness follows from Theorem 2.58 (Sturm).  $\square$

**Complexity analysis:** The complexity of the algorithm is  $O(pq)$  according to the complexity analysis of the Algorithm 8.19 (Signed Remainder Sequence). Indeed, there are only  $O(p)$  extra sign determinations tests to perform.  $\square$



*Remark 9.1.* Note that a much more sophisticated method for computing the Cauchy index is based on ideas from [145, 119]. In this approach, the sign variations in the polynomials of the signed remainder sequence evaluated at  $-\infty$  and  $\infty$  are computed from the quotients and the gcd, with complexity  $O((p+q)\log(p+q)^2) = \tilde{O}(p+q)$  using the fact that the quotients can be computed from the leading terms of the polynomials in the remainder sequence. The same remark applies for Algorithm 9.2.  $\square$

This algorithm gives the following method for computing a Tarski-query. Recall that the Tarski-query of  $Q$  for  $P$  is the number

$$\text{TaQ}(Q, P) = \sum_{x \in \mathbb{R}, P(x)=0} \text{sign}(Q(x)).$$

*Algorithm 9.2. [Remainder Univariate Tarski-query]*

- **Structure:** an ordered field  $K$ .
- **Input:**  $P \in K[X] \setminus \{0\}$ ,  $Q \in K[X]$ .
- **Output:** the Tarski-query  $\text{TaQ}(Q, P)$ .
- **Complexity:**  $O(p(p+q))$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:** Call Algorithm 9.1 (Sylvester Cauchy index) with input  $P$  and  $P'Q$ .

**Proof of correctness:** The correctness follows from Theorem 2.61 (Sylvester's theorem).  $\square$

**Complexity analysis:** Suppose that  $P$  and  $Q$  have respective degree  $p$  and  $q$ . The complexity of the algorithm is  $O(p(p+q))$  according to the complexity analysis of Algorithm 9.1 (Sylvester Cauchy index).  $\square$

**Exercise 9.1.** Design an algorithm computing  $\text{TaQ}(Q, P; a, b)$  with complexity  $O((p+q)^2)$ , where  $\text{TaQ}(Q, P; a, b)$ .

Another algorithm for computing the Cauchy index using the subresultant polynomials is based on Theorem 4.31. Its main advantage is that the bitsize of intermediate computations are much better controlled.

We first compute generalized permanences minus variations (see Notation 4.30).

*Algorithm 9.3. [Generalized Permanences minus Variations]*

- **Structure:** an ordered integral domain  $D$ .
- **Input:**  $s = s_p, \dots, s_0$  be a finite list of elements in  $D$  such that  $s_p \neq 0$ .
- **Output:**  $\text{PmV}(s)$ .
- **Complexity:**  $O(p)$ .
- **Procedure:**
  - Initialize  $n$  to 0.

- Compute the number  $\ell$  of non-zero elements of  $s$  and define the list  $(s'(1), m(1)), \dots, (s'(\ell), m(\ell)) = (s_p, p), (s_q, q), \dots$ , of non-zero elements of  $s$  with their index.
- For every  $i$  from 1 to  $\ell - 1$ , if  $m(i) - m(i + 1)$  is odd

$$n := n + (-1)^{(m(i)-m(i+1))(m(i)-m(i+1)-1)/2} \text{sign}(s'(i) s'(i + 1)).$$

*Algorithm 9.4.* **[Cauchy Index]**

- **Structure:** an ordered integral domain  $D$ .
- **Input:**  $P \in D[X] \setminus \{0\}$ ,  $Q \in D[X]$ .
- **Output:** the Cauchy index  $\text{Ind}((Q/P))$ .
- **Complexity:**  $O(pq)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:** If  $q \geq p$ , replace  $Q$  by the signed pseudo-remainder of  $Q$  and  $P$ . Using Algorithm 8.21 (Signed Subresultant), compute the sequence  $\text{sRes}$  of principal signed subresultant coefficient of  $P$  and  $Q$ , and then compute  $\text{PmV}(\text{sRes}(P, Q))$  (Notation 4.30).

**Proof of correctness:** The correctness follows from Theorem 4.31. □

**Complexity analysis:** The complexity of the algorithm is  $O(pq)$  according to the complexity analysis of Algorithm 8.21 (Signed Subresultant), since there are only  $O(p)$  extra sign evaluations to perform.

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , with coefficients of bitsizes bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $(\tau + \nu)(p + q)$ , where  $\nu$  is the bitsize of  $p + q$ . This follows from Proposition 8.44. □

*Remark 9.2.* Similar ideas to that of [145, 119] (see Remark 9.1) can be used for computing the Cauchy index with complexity  $\tilde{O}(q\tau)$  and binary complexity  $\tilde{O}((p + q)^2\tau)$  [107]. The same remark applies for Algorithm 9.5. □

This algorithm gives the following method for computing Tarski-queries.

*Algorithm 9.5.* **[Univariate Tarski-query]**

- **Structure:** an ordered integral domain  $D$ .
- **Input:**  $P \in D[X] \setminus \{0\}$ ,  $Q \in D[X]$ .
- **Output:** the Tarski-query  $\text{TaQ}(Q, P)$ .
- **Complexity:**  $O((p + q)q)$ , where  $p$  is the degree of  $P$  and  $q$  the degree of  $Q$ .
- **Procedure:**
  - If  $\text{deg}(Q) = 0$ ,  $Q = b_0$ , compute the sequence  $\text{sRes}(P, P')$  of signed subresultant coefficient of  $P$  and  $P'$  using Algorithm 8.21 (Signed Subresultant), and compute  $\text{PmV}(\text{sRes}(P, P'))$  (Definition 4.30). Output

$$\begin{cases} \text{PmV}(\text{sRes}(P, P')) & \text{if } b_0 > 0, \\ -\text{PmV}(\text{sRes}(P, P')) & \text{if } b_0 < 0. \end{cases}$$

- If  $\deg(Q) = 1$ ,  $Q = b_1 X + b_0$ , compute  $R := P' Q - p b_1 P$ , the sequence  $\text{sRes}(P, R)$  of signed subresultant coefficient of  $P$  and  $R$ , using Algorithm 8.21 (Signed Subresultant), and compute  $\text{PmV}(\text{sRes}(P, R))$  (Definition 4.30).
- If  $\deg(Q) > 1$  use Algorithm 8.21 (Signed Subresultant) to compute the sequence  $\text{sRes}(-P' Q, P)$  of signed subresultant coefficient of  $-P' Q$  and  $Q$ , and compute  $\text{PmV}(\text{sRes}(-P' Q, P))$  (Definition 4.30). Output

$$\begin{cases} \text{PmV}(\text{sRes}(-P' Q, P)) + \text{sign}(b_q) & \text{if } q - 1 \text{ is odd,} \\ \text{PmV}(\text{sRes}(-P' Q, P)) & \text{otherwise.} \end{cases}$$

**Proof of correctness:** The correctness follows from Corollary 9.6 and Lemma 9.5. □

**Complexity analysis:** The complexity of the algorithm is  $O((p + q) p)$ , according to the complexity analysis of the Algorithm 8.21 (Signed Subresultant).

Suppose  $P$  and  $Q$  in  $\mathbb{Z}[X]$  with coefficients of bitsizes bounded by  $\tau$ , and denote by  $\nu$  the bitsize of  $2p + q - 1$ .

When  $q > 1$ , the bitsizes of the coefficients of  $P'Q$  are bounded by  $2\tau + \nu$ . When  $q = 1$ , the bitsizes of the coefficients of  $\bar{R}$  are bounded by  $2\tau + 2\nu$ . When  $q = 0$ , the bitsizes of the coefficients of  $P'$  are bounded by  $\tau + \nu$ .

Thus the bitsizes of the integers in the operations performed by the algorithm are bounded by  $(2\tau + 2\nu)(2p + q - 1)$ , according to Proposition 8.44. □

### 9.1.2 Bezoutian and Cauchy Index

We give in this section yet another way of obtaining the Cauchy index. Let  $P$  and  $Q$  be two polynomials with:

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + \dots + a_0 \\ Q &= b_{p-1} X^{p-1} + \dots + b_0, \end{aligned}$$

with  $\deg(P) = p$ ,  $\deg(Q) = q \leq p - 1$ .

**Notation 9.3. [Bezoutian]** The **Bezoutian** of  $P$  and  $Q$  is

$$\text{Bez}(P, Q) = \frac{Q(Y)P(X) - Q(X)P(Y)}{X - Y}.$$

If  $\mathcal{B} = b_1(X), \dots, b_p(X)$  is a basis of  $\mathbb{K}[X]/(P(X))$ ,  $\text{Bez}(P, Q)$  can be uniquely written

$$\text{Bez}(P, Q) = \sum_{i,j=1}^p c_{i,j} b_i(X) b_j(Y).$$

The matrix of  $\text{Bez}(P, Q)$  in the basis  $\mathcal{B}$  is the symmetric matrix with  $i, j$ -th entry the coefficient  $c_{i,j}$  of  $b_i(X)b_j(Y)$  in  $\text{Bez}(P, Q)$ . Note that the signature of the matrix of  $\text{Bez}(P, Q)$  in the basis  $\mathcal{B}$  does not depend of  $\mathcal{B}$  by Sylvester's inertia law (Theorem 4.38). □

**Theorem 9.4.** *The following equalities hold*

$$\begin{aligned} \text{Rank}(\text{Bez}(P, Q)) &= \text{deg}(P) - \text{deg}(\text{gcd}(P, Q)) \\ \text{Sign}(\text{Bez}(P, Q)) &= \text{Ind}(Q/P). \end{aligned}$$

The proof of the Theorem will use the following results.

**Lemma 9.5.** *Suppose  $s = s_0, \dots, s_{c-2}, s_{c-1}, \dots, s_{2n-2}$ , with  $2n - 1 \geq c \geq n$ , and  $s_0 = \dots, s_{c-2} = 0, s_{c-1} \neq 0$ , and let  $H$  be the  $n \times n$  matrix defined by  $h_{i,j} = s_{i+j-2}$ . Then*

$$\begin{aligned} \text{Rank}(H) &= 2n - c, \\ \text{Sign}(H) &= \begin{cases} \text{sign}(s_{c-1}) & \text{if } c \text{ is odd,} \\ 0 & \text{if } c \text{ is even.} \end{cases} \end{aligned}$$

The proof of the lemma is based on the following proposition.

**Proposition 9.6.** *Let  $H$  be a semi-algebraic continuous mapping from an interval  $I$  of  $\mathbb{R}$  into the set of symmetric matrix of dimension  $n$ . If, for every  $t \in I$ , the rank of  $H(t)$  is always equal to the same value, then, for every  $t \in I$ , the signature of  $H(t)$  is always equal to the same value.*

**Proof:** Let  $r$  be the rank of  $H(t)$ , for every  $t \in I$ . The number of zero eigenvalues of  $H(t)$  is  $n - r$  for every  $t \in I$ , by Corollary 4.44. The number of positive and negative eigenvalues of  $H(t)$  is thus also constant, since roots vary continuously (see Theorem 5.12 (Continuity of roots)). Thus, by Corollary 4.44, for every  $t \in I$ , the signature of  $H(t)$  is always equal to the same value. □

**Proof of Lemma 9.5:** Let  $s_c = \dots = s_{2n-2} = 0$ . In this special case, the rank of  $H$  is obviously  $2n - c$ . Since the associated quadratic form is

$$\Phi = \sum_{i=c+1-n}^n s_{c-1} f_i f_{c+1-i}$$

and, if  $c + 1 - i \neq i$ ,

$$4f_i f_{c+1-i} = (f_i + f_{c+1-i})^2 - (f_i - f_{c+1-i})^2,$$

it is easy to see that the signature of  $\Phi$  is 0 if  $c$  is even, and is 1 (resp.  $-1$ ) if  $s_{c-1} > 0$  (resp.  $s_{c-1} < 0$ ).

Defining, for  $t \in [0, 1]$ ,  $s_t = 0, \dots, 0, s_{c-1}, t s_c, \dots, t s_{2n-2}$  the quadratic form with associated matrix  $H_t$  defined by  $h_{t,i,j} = s_{t,i+j-2}$  is of rank  $2n - c$  for every  $t \in [0, 1]$ , since the  $c - n$  first columns of  $H_t$  are zero and its  $2n - c$  last columns are clearly independent. Thus the rank of  $H_t$  is constant as  $t$  varies. Thus by Proposition 9.6 the signature of  $H_t$  is constant as  $t$  varies. This proves the claim, taking  $t = 1$ .  $\square$

We denote by  $R$  and  $C$  the remainder and quotient of the euclidean division of  $P$  by  $Q$ .

**Lemma 9.7.** *The following equalities hold.*

$$\begin{aligned} \text{Rank}(\text{Bez}(P, Q)) &= \text{Rank}(\text{Bez}(Q, -R)) + p - q \\ \text{Sign}(\text{Bez}(P, Q)) &= \begin{cases} \text{Sign}(\text{Bez}(Q, -R)) + \text{sign}(a_p b_q) & \text{if } p - q \text{ is odd,} \\ \text{Sign}(\text{Bez}(Q, -R)) & \text{if } p - q \text{ is even.} \end{cases} \end{aligned}$$

**Proof:** We consider the matrix  $M(P, Q)$  of coefficients of  $\text{Bez}(P, Q)$  in the canonical basis

$$X^{p-1}, \dots, 1,$$

and the matrix  $M'(P, Q)$  of coefficients of  $\text{Bez}(P, Q)$  in the basis

$$X^{p-q-1} Q(X), \dots, X Q(X), Q(X), X^{q-1}, \dots, 1.$$

Let  $c = p - q$ ,  $C = u_c X^c + \dots + u_0$ ,  $s = \overbrace{0, \dots, 0}^{c-1 \text{ times}}, u_c, \dots, u_1$  of length  $2c - 1$  and  $H$  the  $c \times c$  matrix defined by  $h_{i,j} = s_{i+j-2}$ . Since  $P = CQ + R$ , and  $\text{deg}(R) < q$ ,

$$\text{Bez}(P, Q) = \frac{C(X) - C(Y)}{X - Y} Q(Y) Q(X) + \text{Bez}(Q, -R),$$

the matrix  $M'(P, Q)$  is the block matrix

$$\begin{bmatrix} H & 0 \\ 0 & M(Q, -R) \end{bmatrix}.$$

The claim follows from Lemma 9.5 since the leading coefficient of  $C$  is  $a_p/b_q$ .  $\square$

**Proof of Theorem 9.4:** The proof of the Theorem proceeds by induction on the number  $n$  of elements with distinct degrees in the signed subresultant sequence.

If  $n = 2$ ,  $Q$  divides  $P$  and  $R = 0$ . We have

$$\text{Rank}(\text{Bez}(P, Q)) = \text{deg}(P) - \text{deg}(Q).$$

We also have

$$\text{Ind}(Q/P) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p - q \text{ is odd,} \\ 0 & \text{if } p - q \text{ is even.} \end{cases}$$

by Lemma 4.34 and

$$\text{Sign}(\text{Bez}(P, Q)) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p - q \text{ is odd,} \\ 0 & \text{if } p - q \text{ is even.} \end{cases}$$

by Lemma 9.7.

Let us suppose that the Theorem holds for  $n - 1$  and consider  $P$  and  $Q$  such that their signed subresultant sequence has  $n$  elements with distinct degrees. The signed subresultant sequence of  $Q$  and  $-R$  has  $n - 1$  elements with distinct degrees and by induction hypothesis,

$$\begin{aligned} \text{Rank}(\text{Bez}(Q, -R)) &= \text{deg}(Q) - \text{deg}(\text{gcd}(Q, -R)) \\ \text{Sign}(\text{Bez}(Q, -R)) &= \text{Ind}(-R/Q) \end{aligned}$$

By Lemma 9.4 and Lemma 9.7, since  $\text{gcd}(P, Q) = \text{gcd}(Q, -R)$ .

$$\begin{aligned} \text{Rank}(\text{Bez}(P, Q)) &= \text{deg}(P) - \text{deg}(\text{gcd}(P, Q)) \\ \text{Sign}(\text{Bez}(Q, -R)) &= \text{Ind}(Q/P) \end{aligned}$$

□

Theorem 9.4 has the following corollaries.

**Corollary 9.8.** *Let  $P$  and  $Q$  be polynomials in  $D[X]$  and  $R$  the remainder of  $P'Q$  divided by  $P$ . Then  $\text{Sign}(\text{Bez}(P, R)) = \text{TaQ}(Q, P)$ .*

**Proof:** Apply Theorem 9.4 and Proposition 2.57, noticing that

$$\text{Ind}(P'Q/P) = \text{Ind}(R/P). \quad \square$$

**Corollary 9.9.** *Let  $P$  be a polynomial in  $D[X]$ . Then  $\text{Sign}(\text{Bez}(P, P'))$  is the number of roots of  $P$  in  $R$ .*

It follows immediately from Theorem 9.4 that the determinant of the matrix of  $\text{Bez}(P, Q)$  in the canonical basis  $X^{p-1}, \dots, 1$ , is 0 if and only if  $\text{deg}(\text{gcd}(P, Q)) > 0$ . On the other hand, by Proposition 4.15  $\text{Res}(P, Q) = 0$  if and only if  $\text{deg}(\text{gcd}(P, Q)) > 0$ . This suggests a close connection between the determinant of the matrix of  $\text{Bez}(P, Q)$  in the canonical basis  $X^{p-1}, \dots, 1$ , and  $\text{Res}(P, Q)$ .

**Proposition 9.10.** *Let  $M(P, Q)$  be the matrix of coefficients of  $\text{Bez}(P, Q)$  in the canonical basis  $X^{p-1}, \dots, 1$ . Then*

$$\det(M(P, Q)) = \varepsilon_p a_p^{p-q} \text{Res}(P, Q),$$

with  $\varepsilon_i = (-1)^{i(i-1)/2}$ .

**Proof:** If  $\text{deg}(\text{gcd}(P, Q)) > 0$ , both quantities are 0, so the claim is true.

Suppose now that  $\text{gcd}(P, Q)$  is a constant. The proof is by induction on the number  $n$  of elements in the signed remainder sequence.

If  $n = 2$ , and  $Q = b, q = 0$ ,

$$\begin{aligned} \text{Res}(P, Q) &= b^{p-q}, \\ \det(M(P, Q)) &= \varepsilon_p a_p^p b^p \end{aligned}$$

and the claim holds.

Suppose that the claim holds for  $n - 1$ . According to Proposition 8.35 and Equation (4.3) in Notation 4.26

$$\varepsilon_p \text{Res}(P, Q) = \varepsilon_q \varepsilon_{p-q} b_q^{p-q} \text{Res}(Q, -R). \tag{9.1}$$

On the other hand, by the proof of Lemma 9.7, and using its notation, the matrix  $M'(P, Q)$  of coefficients of  $\text{Bez}(P, Q)$  in the basis

$$X^{p-q-1} Q, \dots, X Q, Q, X^{q-1}, \dots, 1$$

is the matrix

$$\begin{bmatrix} H & 0 \\ 0 & M(Q, -R) \end{bmatrix}$$

with, for  $1 \leq i \leq p - q$ ,  $h_{p-q+1-i, i} = a_p/b_q$ . Thus, using the fact that the leading coefficient of  $Q$  is  $b_q$ ,

$$\det(\text{Bez}(P, Q)) = \varepsilon_{p-q} a_p^{p-q} b_q^{p-q} \det(\text{Bez}(Q, -R)). \tag{9.2}$$

The induction hypothesis

$$\det(M(Q, -R)) = \varepsilon_q b_q^{q-r} \text{Res}(Q, -R),$$

Equation (9.1) and Equation (9.2), imply the claim for  $P, Q$ . □

### 9.1.3 Signed Subresultant Sequence and Cauchy Index on an Interval

We show that the Cauchy index on an interval can be expressed in terms of appropriately counted sign variations in the signed subresultant sequence. The next definitions introduce the sign counting function to be used.

**Notation 9.11.** Let  $s = s_n, 0, \dots, 0, s'$ , be a finite sequence of elements in an ordered field  $K$  such that  $s_n \neq 0$ ,  $s' = \emptyset$  or  $s' = s_m, \dots, s_0, s_m \neq 0$ . The **modified number of sign variations** in  $s$  is defined inductively as follows

$$\text{MVar}(s) = \begin{cases} 0 & \text{if } s' = \emptyset, \\ \text{MVar}(s') + 1 & \text{if } s_n s_m < 0, \\ \text{MVar}(s') + 2 & \text{if } s_n s_m > 0 \text{ and } n - m = 3 \\ \text{MVar}(s') & \text{if } s_n s_m > 0 \text{ and } n - m \neq 3, \end{cases}$$

In other words, we modify the usual definition of the number of sign variations by counting 2 sign variations for the groups:  $+, 0, 0, +$  and  $-, 0, 0, -$ .

Let  $\mathcal{P} = P_0, P_1, \dots, P_d$  be a sequence of polynomials in  $D[X]$  and  $a$  be an element of  $R \cup \{-\infty, +\infty\}$  which is not a root of  $\gcd(\mathcal{P})$ . Then  $MVar(\mathcal{P}; a)$ , the **modified number of sign variations** of  $\mathcal{P}$  at  $a$ , is the number defined as follows:

- Delete from  $\mathcal{P}$  those polynomials that are identically 0 to obtain the sequence of polynomials  $\mathcal{Q} = Q_0, \dots, Q_s$  in  $D[X]$ ,
- Define  $MVar(\mathcal{P}; a)$  as  $MVar(Q_0(a), \dots, Q_s(a))$ .

Let  $a$  and  $b$  be elements of  $R \cup \{-\infty, +\infty\}$  which are not roots of  $\gcd(\mathcal{P})$ . The difference between the number of modified sign variations in  $\mathcal{P}$  at  $a$  and  $b$  is denoted by

$$MVar(\mathcal{P}; a, b) = MVar(\mathcal{P}; a) - MVar(\mathcal{P}; b). \quad \square$$

For example, if  $\mathcal{P} = X^5, X^2 - 1, 0, X^2 - 1, X + 2, 1$ , the modified number of sign variations of  $\mathcal{P}$  at 1 is 2 while the number of signs variations of  $\mathcal{P}$  at 1 is 0.

**Theorem 9.12.**

$$MVar(\text{sResP}(P, Q); a, b) = \text{Ind}(Q/P; a, b).$$

Note that when polynomials of all possible degrees  $\leq p$  appear in the remainder sequence, Theorem 9.12 is an immediate consequence of Theorem 2.58, since the signed remainder sequence and the signed subresultant sequence are proportional up to squares by Corollary 8.33.

The proof of Theorem 9.12 uses the following lemma.

**Lemma 9.13.** *Let  $R = \text{Rem}(P, Q)$  and let  $\sigma(a)$  be the sign of  $PQ$  at  $a$  and  $\sigma(b)$  be the sign of  $PQ$  at  $b$ . Then*

$$\begin{aligned} & MVar(\text{sResP}(P, Q); a, b) \\ = & \begin{cases} MVar(\text{sResP}(Q, -R); a, b) + \sigma(b) & \text{if } \sigma(a)\sigma(b) = -1, \\ MVar(\text{sResP}(Q, -R); a, b) & \text{if } \sigma(a)\sigma(b) = 1. \end{cases} \end{aligned}$$

**Proof:** We denote  $L = \text{sResP}(P, Q)$  and  $L' = \text{sResP}(Q, -R)$ .

Suppose that  $x$  be not a root of  $P, Q$ , or  $R$ . According to Proposition 8.35, and the conventions in Notation 8.29,

$$\begin{aligned} \text{sResP}_p(P, Q) &= P, \\ \text{sResP}_{p-1}(P, Q) &= Q, \\ \text{sResP}_q(P, Q) &= \varepsilon_{p-q} b_q^{p-q-1} Q, \\ \text{sResP}_q(Q, -R) &= Q, \\ \text{sResP}_{q-1}(P, Q) &= -\varepsilon_{p-q} b_q^{p-q+1} R, \\ \text{sResP}_{q-1}(Q, -R) &= -R, \\ \text{sResP}_j(P, Q) &= \varepsilon_{p-q} b_q^{p-r} \text{sResP}_j(Q, -R), \quad j \leq r. \end{aligned}$$

Hence for every  $x$  which is not a root of  $P, Q$  and  $-R$ , in particular for  $a$  and  $b$ , the following holds, denoting by  $c_r$  the leading coefficient of  $-R$ .



If  $P(x)Q(x) > 0$

– if  $\varepsilon_{q-r}c_r^{q-r-1} > 0$

$$\text{MVar}(L; x) = \begin{cases} \text{MVar}(L'; x) + 2 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} < 0 \text{ and } b_q^{q-r-1} < 0, \\ \text{MVar}(L'; x) + 1 & \text{if } \varepsilon_{p-q}b_q^{p-r} < 0, \\ \text{MVar}(L'; x) & \text{if } \varepsilon_{p-q}b_q^{p-q-1} > 0 \text{ and } b_q^{q-r-1} > 0, \end{cases}$$

– if  $\varepsilon_{q-r}c_r^{q-r-1} < 0$

$$\text{MVar}(L; x) = \begin{cases} \text{MVar}(L'; x) + 1 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} < 0 \text{ and } b_q^{q-r-1} > 0, \\ \text{MVar}(L'; x) & \text{if } \varepsilon_{p-q}b_q^{p-r} > 0, \\ \text{MVar}(L'; x) - 1 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} > 0 \text{ and } b_q^{q-r-1} < 0. \end{cases}$$

If  $P(x)Q(x) < 0$ ,

– if  $\varepsilon_{q-r}c_r^{q-r-1} > 0$

$$\text{MVar}(L; x) = \begin{cases} \text{MVar}(L'; x) + 3 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} < 0 \text{ and } b_q^{q-r-1} < 0, \\ \text{MVar}(L'; x) + 2 & \text{if } \varepsilon_{p-q}b_q^{p-r} < 0, \\ \text{MVar}(L'; x) + 1 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} > 0 \text{ and } b_q^{q-r-1} > 0, \end{cases}$$

– if  $\varepsilon_{q-r}c_r^{q-r-1} < 0$

$$\text{MVar}(L; x) = \begin{cases} \text{MVar}(L'; x) + 2 & \text{if } \varepsilon_{p-q}b_q^{p-q-1} < 0 \text{ and } b_q^{q-r-1} > 0, \\ \text{MVar}(L'; x) + 1 & \text{if } \varepsilon_{p-q}b_q^{p-r} > 0, \\ \text{MVar}(L'; x) & \text{if } \varepsilon_{p-q}b_q^{p-q-1} > 0 \text{ and } b_q^{q-r-1} < 0 \end{cases}$$

The lemma follows easily. □

**Proof of Theorem 9.12:** We can assume without loss of generality that  $a$  and  $b$  are not roots of a non-zero polynomial in the signed subresultant sequence. Indeed if  $a < a' < b' < b$  with  $(a, a']$  and  $[b', b)$  containing no root of the polynomials in the signed subresultant sequence,

$$\text{Ind}(Q/P; a, b) = \text{Ind}(Q/P; a', b').$$

We also have

$$\text{MVar}(\text{sResP}(P, Q); a, b) = \text{MVar}(\text{sResP}(P, Q); a', b').$$

Indeed if  $a$  is a root of of  $\text{sResP}_{j-1}$ ,

– when  $\text{sResP}_{j-1}$  is non-defective, we have

$$\begin{aligned} \text{MVar}(\text{sResP}_{j-2}, \text{sResP}_{j-1}, \text{sResP}_j; a) &= 1, \\ \text{MVar}(\text{sResP}_{j-2}, \text{sResP}_{j-1}, \text{sResP}_j; a') &= 1, \end{aligned}$$

– when  $\text{sResP}_{j-1}(P, Q)$  is defective of degree  $k$ , we have

$$\begin{aligned} & \text{MVar}(\text{sResP}_j, \text{sResP}_{j-1}, \text{sResP}_k, \text{sResP}_{k-1}; a') \\ &= \text{MVar}(\text{sResP}_j, \text{sResP}_{j-1}, \text{sResP}_k, \text{sResP}_{k-1}; a) \\ &= \begin{cases} 2 & \text{if } \text{sResP}_j(a)\text{sResP}_{k-1}(a) > 0 \\ 1 & \text{if } \text{sResP}_j(a)\text{sResP}_{k-1}(a) < 0. \end{cases} \end{aligned}$$

The proof of the theorem proceeds by induction on the number  $n$  of elements with distinct degrees in the signed subresultant sequence. The base case  $n = 2$  corresponds to  $\deg(Q) = 0$ ,  $R = 0$  and follows from Lemma 9.13 and Lemma 2.60. Let us suppose that the Theorem holds for  $n - 1$  and consider  $P$  and  $Q$  such that their signed subresultant sequence has  $n$  elements with distinct degrees. The signed subresultant sequence of  $Q$  and  $-R$  has  $n - 1$  elements with distinct degrees and by the induction hypothesis,

$$\text{MVar}(\text{sResP}(Q, -R); a, b) = \text{Ind}(-R/Q; a, b).$$

So, by Lemma 9.13 and Lemma 2.60,

$$\text{MVar}(\text{sResP}(P, Q); a, b) = \text{Ind}(Q/P; a, b). \quad \square$$

**Corollary 9.14.** *Let  $P, Q \in \mathbb{D}[X]$ . Let  $R$  be the remainder of  $P'Q$  and  $P$ . If  $a < b$  are elements of  $\mathbb{R} \cup \{-\infty, +\infty\}$  that are not roots of  $P$ , then*

$$\text{MVar}(\text{sResP}(P, R); a, b) = \text{TaQ}(Q, P; a, b).$$

**Proof:** Apply Theorem 9.12 and Proposition 2.57, since

$$\text{Ind}(P'Q/P; a, b) = \text{Ind}(R/P; a, b),$$

by Remark 2.55. □

**Corollary 9.15.** *Let  $P$  be a polynomial in  $\mathbb{D}[X]$ . If  $a < b$  are elements of  $\mathbb{R} \cup \{-\infty, +\infty\}$  which are not roots of  $P$ , then  $\text{MVar}(\text{sResP}(P, P'); a, b)$  is the number of roots of  $P$  in  $(a, b)$ .*

**Exercise 9.2.** Using Algorithm 8.21 (Signed Subresultant), design an algorithm computing  $\text{TaQ}(Q, P; a, b)$  with complexity  $O((p + q)p)$ . If  $P \in \mathbb{Z}[X]$ ,  $Q \in \mathbb{Z}[X]$ ,  $a, b$  are rational numbers, and  $\tau$  is a bound on the bitsize of the coefficients of  $P$  and  $Q$  and on  $a$  and  $b$ , estimate the bitsize of the rationals computed by this algorithm.

## 9.2 Hankel Matrices

Hankel matrices are important because of their relation with rational functions and sequences satisfying linear recurrence relations. We define Hankel matrices and quadratic forms and indicate how to compute the corresponding signature.

### 9.2.1 Hankel Matrices and Rational Functions

Hankel matrices are symmetric matrix with equal entries on the anti-diagonals. More precisely **Hankel matrices** of size  $p$  are matrices with entries  $a_{i+1,j+1}$  ( $i$  from 0 to  $p-1$ ,  $j$  from 0 to  $p-1$ ) such that  $a_{i+1,j+1} = a_{i'+1,j'+1}$  whenever  $i+j = i'+j'$ .

A typical Hankel matrix is the matrix  $\text{Newt}_k(P)$  considered in Chapter 4.

**Notation 9.16. [Hankel]** Let  $\bar{s} = s_0, \dots, s_n, \dots$  be an infinite sequence. We denote  $\bar{s}_n = s_0, \dots, s_{2n-2}$ , and  $\text{Han}(\bar{s}_n)$  the Hankel matrix whose  $i+1, j+1$  entry is  $s_{i+j}$  for  $0 \leq i, j \leq n-1$ , and by  $\text{han}(\bar{s}_n)$  the determinant of  $\text{Han}(\bar{s}_n)$ .  $\square$

**Theorem 9.17.** *Let  $K$  be a field. Let  $\bar{s} = s_0, \dots, s_n, \dots$  be an infinite sequence of elements of  $K$  and  $p \in \mathbb{N}$ . The following conditions are equivalent:*

- a) *The elements  $s_0, \dots, s_n, \dots$  satisfy a linear recurrence relation of order  $p$  with coefficients in  $K$*

$$a_p s_n = -a_{p-1} s_{n-1} - \dots - a_0 s_{n-p}, \quad (9.3)$$

$$a_p \neq 0, n \geq p.$$

- b) *There exists a polynomial  $P \in K[X]$  of degree  $p$  and a linear form  $\lambda$  on  $K[X]/(P)$  such that  $\lambda(X^i) = s_i$  for every  $i \geq 0$ .*  
 c) *There exist polynomials  $P, Q \in K[X]$  with  $\deg(Q) < \deg(P) = p$  such that*

$$Q/P = \sum_{j=0}^{\infty} s_j / X^{j+1} \quad (9.4)$$

- d) *There exists an  $r \leq p$  such that the ranks of all the Hankel matrices  $\text{Han}(\bar{s}_r), \text{Han}(\bar{s}_{r+1}), \text{Han}(\bar{s}_{r+2}), \dots$  are equal to  $r$ .*  
 e) *There exists an  $r \leq p$  such that*

$$\text{han}(\bar{s}_r) \neq 0, \quad \forall n > r \text{ han}(\bar{s}_n) = 0.$$

A sequence satisfying the equivalent properties of Theorem 9.17 is a **linear recurrent sequence of order  $p$** .

The proof of the theorem uses the following definitions and results. Let

$$P = a_p X^p + a_{p-1} X^{p-1} + \dots + a_1 X + a_0$$

be a polynomial of degree  $p$ . The Horner polynomials associated to  $P$  are defined inductively by

$$\begin{aligned} \text{Hor}_0(P, X) &= a_p, \\ &\vdots \\ \text{Hor}_i(P, X) &= X \text{Hor}_{i-1}(P, X) + a_{p-i}, \\ &\vdots \end{aligned}$$

for  $i = 0, \dots, p-1$  (see Notation 8.6).

The Horner polynomials

$$\text{Hor}_0(P, X), \dots, \text{Hor}_{p-1}(P, X)$$

are obviously a basis of the quotient ring  $\mathbb{K}[X]/(P)$ .

The **Kronecker form** is the linear form  $\ell_P$  defined on  $\mathbb{K}[X]/(P(X))$ , by

$$\ell_P(1) = \dots = \ell_P(X^{p-2}) = 0, \ell_P(X^{p-1}) = 1/a_p.$$

If  $Q \in \mathbb{K}[X]$ ,  $\text{Rem}(Q, P)$  is the canonical representative of its equivalence class in  $\mathbb{K}[X]/(P(X))$ , and  $\ell_P(Q)$  denotes  $\ell_P(\text{Rem}(Q, P))$ .

**Proposition 9.18.** For  $0 \leq i \leq p-1, 0 \leq j \leq p-1$ ,

$$\ell_P(X^j \text{Hor}_{p-1-i}(P, X)) = \begin{cases} 1, & j=i, \\ 0, & j \neq i. \end{cases}$$

**Proof:** The claim is clear from the definitions if  $j \leq i$ .

If  $i < j \leq p-1$ , since

$$a_p X^p + \dots + a_{i+1} X^{i+1} = -(a_i X^i + \dots + a_0 \text{ mod } P(X)),$$

we have

$$\begin{aligned} X^{i+1} \text{Hor}_{p-1-i}(P, X) &= -(a_i X^i + \dots + a_0) \text{ mod } P(X) \\ X^j \text{Hor}_{p-1-i}(P, X) &= -X^{j-i-1} (a_i X^i + \dots + a_0) \text{ mod } P(X), \end{aligned}$$

and, by definition of  $\ell_P$

$$\ell_P(X^j \text{Hor}_{p-1-i}(P, X)) = -\ell_P(X^{j-i-1} (a_i X^i + \dots + a_0)) = 0. \quad \square$$

**Corollary 9.19.** For every  $Q \in \mathbb{K}[X]$ ,

$$Q = \sum_{i=0}^{p-1} \ell_P(Q X^i) \text{Hor}_{p-1-i}(P, X) \text{ mod } P(X). \tag{9.5}$$

**Proof:** By (9.18),

$$\text{Hor}_{p-1-j}(P, X) = \sum_{i=0}^{p-1} \ell_P(\text{Hor}_{p-1-j}(P, X) X^i) \text{Hor}_{p-1-i}(P, X) \text{ mod } P(X).$$

The claim follows by the linearity of  $\ell_P$  after expressing  $Q_1 = \text{Rem}(Q, P)$  in the Horner basis. □

**Proof of Theorem 9.17:**  $a) \Rightarrow c)$ : Take

$$\begin{aligned} P &= a_p X^p + a_{p-1} X^{p-1} + \dots + a_0, \\ Q &= s_0 \text{Hor}_{p-1}(P, X) + \dots + s_{p-1} \text{Hor}_0(P, X). \end{aligned}$$

Note that if  $Q = b_{p-1}X^{p-1} + \dots + b_0$ , then

$$b_{p-n-1} = a_p s_n + \dots + a_{p-n} s_0 \tag{9.6}$$

for  $0 \leq n \leq p-1$ , identifying the coefficients of  $X^{p-n-1}$  on both sides of (9.5). Let  $t_n$  be the infinite sequence defined by the development of the rational fraction  $Q/P$  as a series in  $1/X$ :

$$Q/P = \left( \sum_{n=0}^{\infty} t_n/X^{n+1} \right). \tag{9.7}$$

Thus,

$$Q = \left( \sum_{n=0}^{\infty} t_n/X^{n+1} \right) P. \tag{9.8}$$

Identifying the coefficients of  $X^{p-n-1}$  on both sides of (9.8) proves that for  $0 \leq n < p$

$$b_{p-n-1} = a_p t_n + \dots + a_{p-n} t_0,$$

and for  $n \geq p$

$$a_p t_n + a_{p-1} t_{n-1} + \dots + a_0 t_{n-p} = 0.$$

Since  $a_p \neq 0$ , the sequences  $s_n$  and  $t_n$  have the same  $p$  initial values and satisfy the same recurrence relation. Hence they coincide.

$c) \Rightarrow b)$ : For  $i = 0, \dots, p-1$ , take  $\lambda(X^i) = \ell_P(Q X^i)$ , where  $\ell_P$  is the Kronecker form. Since  $Q = \sum_{k=0}^{p-1} s_k \text{Hor}_{p-1-k}(P, X)$ , using (9.5),

$$\lambda(X^j) = \ell_P(Q X^j) = \sum_{k=0}^{p-1} s_k \ell_P(X^j \text{Hor}_{p-1-k}(P, X)) = s_j. \tag{9.9}$$

$b) \Rightarrow a)$  is clear, taking for  $a_i$  the coefficient of  $X^i$  in  $P$  and noticing that

$$\begin{aligned} a_p s_n &= \lambda(a_p X^n) \\ &= -\lambda(a_{p-1} X^{n-1} + \dots + a_0 X^{n-p}) \\ &= -a_{p-1} \lambda(X^{n-1}) + \dots + a_0 \lambda(X^{n-p}) \\ &= -a_{p-1} s_{n-1} + \dots + a_0 s_{n-p}. \end{aligned}$$

$a) \Rightarrow d)$ : For  $(n, m) \in \mathbb{N}^2$ , define  $v_{m,n}$  as the vector  $(s_m, \dots, s_{m+n})$ . The recurrence relation (9.3) proves that for  $m \geq p$ ,

$$a_p v_{m,n} = -a_{p-1} v_{m-1,n} - \dots - a_0 v_{m-p,n}.$$

It is easy to prove by induction on  $n$  that the vector space generated by  $v_{0,n}, \dots, v_{n,n}$  is of dimension  $\leq p$ , which proves the claim.

$d) \Rightarrow e)$ : is clear.

$e) \Rightarrow a)$ : Let  $r \leq p$  be such that

$$\text{han}(\bar{s}_r) \neq 0, \forall n > r \text{ han}(\bar{s}_n) = 0.$$

Then the vector  $v_{n-r,r}$ ,  $n \geq r$ , is a linear combination of  $v_{0,r}, \dots, v_{r-1,r}$ . Developing the determinant of the square matrix with columns

$$v_{0,r}, \dots, v_{r-1,r}, v_{n-r,r}$$

on the last columns gives

$$\mu_r s_n + \mu_{r-1} s_{n-1} + \dots + \mu_0 s_{n-r} = 0,$$

with  $\mu_i$  the cofactor of the  $i - 1$ -th element of the last column. Since  $\mu_r \neq 0$ , take  $a_i = \mu_{p-r+i}$ . □

### 9.2.2 Signature of Hankel Quadratic Forms

**Hankel quadratic forms** are quadratic forms associated to Hankel matrices. We design an algorithm for computing the signature of a general Hankel form. Note that we have already seen a special case where the signature of a Hankel quadratic form is particularly easy to determine in Lemma 9.5.

Given  $\bar{s}_n = s_0, \dots, s_{2n-2}$ , we write

$$\overline{\text{Han}}(\bar{s}_n) = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} s_{i+j} f_i f_j.$$

Note that the Hermite quadratic form seen in Chapter 4 is a Hankel form. Let

$$Q/P = \sum_{j=0}^{\infty} s_j / X^{j+1} \in K[[1/X]],$$

and  $\text{Hor}_{p-1}(P, X), \dots, \text{Hor}_0(P, X)$  the Horner basis of  $P$ . We indicate now the relationship between the Hankel quadratic form

$$\overline{\text{Han}}(\bar{s}_p) = \sum_{i=0}^{p-1} \sum_{j=0}^{p-1} s_{i+j} f_i f_j.$$

and the quadratic form  $\text{Bez}(P, Q)$  (see Notation 9.3).

**Proposition 9.20.** *The matrix of coefficients of  $\text{Bez}(P, Q)$  in the Horner basis  $\text{Hor}_{p-1}(P, X), \dots, \text{Hor}_0(P, X)$  is  $\overline{\text{Han}}(\bar{s}_p)$ , i.e.*

$$\text{Bez}(P, Q) = \sum_{i,j=0}^{p-1} s_{i+j} \text{Hor}_{p-1-i}(P, X) \text{Hor}_{p-1-j}(P, Y). \tag{9.10}$$

**Proof:** Indeed,

$$\text{Bez}(P, Q) = \frac{Q(Y) - Q(X)}{X - Y} P(X) + Q(X) \frac{P(X) - P(Y)}{X - Y} \text{mod } P(X),$$

which implies

$$\text{Bez}(P, Q) = Q(X) \frac{P(X) - P(Y)}{X - Y} \text{mod } P(X),$$

noting that

$$\frac{P(X) - P(Y)}{X - Y} = \sum_{i=0}^{p-1} X^i \text{Hor}_{p-1-i}(P, Y),$$

using Corollary 9.19 and Equation (9.9),

$$\begin{aligned} & Q(X) \frac{P(X) - P(Y)}{X - Y} \\ &= \sum_{i=0}^{p-1} Q(X) X^i \text{Hor}_{p-1-i}(P, Y) \\ &= \sum_{i=0}^{p-1} \sum_{j=0}^{p-1} \ell_P(Q(X) X^{i+j}) \text{Hor}_{p-1-i}(P, X) \text{Hor}_{p-1-j}(P, Y) \text{ mod } P(X) \\ &= \sum_{i=0}^{p-1} \sum_{j=0}^{p-1} s_{i+j} \text{Hor}_{p-1-i}(P, X) \text{Hor}_{p-1-i}(P, Y) \text{ mod } P(X). \end{aligned}$$

This proves (9.10). □

*Remark 9.21.* So, by Proposition 4.55, the Hermitic quadratic form  $\text{Her}(P, Q)$  is nothing but  $\text{Bez}(P, R)$  expressed in the Horner basis of  $P$ , with  $R$  the remainder of  $P'Q$  by  $P$ . This proves that Theorem 9.4 is a generalization of Theorem 4.57. □

Let  $\bar{s}_n = 0, \dots, 0, s_{c-1}, \dots, s_{2n-2}, s_{c-1} \neq 0, c < n$ , and define the series  $S$  in  $1/X$  by

$$S = \sum_{j=0}^{2n-2} s_j / X^{j+1}.$$

Consider the inverse  $S^{-1}$  of  $S$ , which is a Laurent series in  $1/X$  and define  $C \in \mathbb{K}[X], T \in \mathbb{K}[[1/X]]$  by  $S^{-1} = C + T$ , i.e.  $(C + T)S = 1$ . Since  $S$  starts with  $s_{c-1}/X^c$ , it is clear that  $\text{deg}(C) = c$ . Let

$$C = u_c X^c + \dots + u_0,$$

and  $\bar{u} = \overbrace{0, \dots, 0}^{c-1 \text{ times}}, u_c, \dots, u_1$  of length  $2c - 1$ , and  $T = \sum_{j=0}^{\infty} t_j / X^{j+1}$ . Note that  $u_c = 1/s_{c-1}$ . Let  $\bar{t}_{n-c} = t_0, \dots, t_{2n-2c-2}$ .

**Lemma 9.22.**

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \begin{cases} \text{Sign}(\text{Han}(\bar{t}_{n-c})) + \text{sign}(s_{c-1}) & \text{if } c \text{ is odd,} \\ \text{Sign}(\text{Han}(\bar{t}_{n-c})) & \text{if } c \text{ is even.} \end{cases}$$

Lemma 9.22 is a consequence of the following Lemma, which uses Toeplitz matrices, i.e. matrices with equal entries on the parallels to the diagonal. More precisely a **Toeplitz matrix** of size  $n$  is a matrix with entries  $a_{i,j}$  ( $i$  from 1 to  $n$ ,  $j$  from 1 to  $n$ ) such that  $a_{i,j} = a_{i',j'}$  whenever  $j - i = j' - i'$ .

**Notation 9.23. [Toeplitz]** Let  $\bar{v} = v_0, v_1, \dots, v_{n-1}$ . We denote by  $\text{To}(\bar{v})$  the triangular Toeplitz matrix of size  $n$  whose  $i, j$ -th entry is  $v_{j-i}$  for  $0 \leq i, j \leq n$ , with  $j - i \geq 0$ , 0 otherwise. □

**Lemma 9.24.**

$$\text{Han}(\bar{s}_n) = \text{To}(\bar{v})^t \begin{bmatrix} \text{Han}(\bar{u}) & 0 \\ 0 & \text{Han}(\bar{t}_{n-c}) \end{bmatrix} \text{To}(\bar{v}),$$

with  $\bar{v} = s_{c-1}, \dots, s_{n+c-2}$ .

We first explain how Lemma 9.22 is a consequence of Lemma 9.24 before proving the lemma itself.

**Proof of Lemma 9.22:** Using Lemma 9.24, the quadratic forms associated to  $\text{Han}(\bar{s}_n)$  and

$$\begin{bmatrix} \text{Han}(\bar{u}) & 0 \\ 0 & \text{Han}(\bar{t}_{n-c}) \end{bmatrix}$$

have the same signature, and

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{Sign}(\overline{\text{Han}}(\bar{t}_{n-c})) + \text{Sign}(\overline{\text{Han}}(\bar{u})).$$

The claim follows from Lemma 9.5, since, as noted above,  $u_c = 1/s_{c-1}$ . □

The proof of Lemma 9.24 requires some preliminary work.

Let  $P = a_p X^p + \dots + a_0$  and  $Q = b_q X^q + \dots + b_0$ ,  $q = \deg(Q) < p = \deg(P)$ , such that

$$Q/P = \sum_{j=0}^{\infty} s_j / X^{j+1} \in K[[1/X]].$$

Then,  $s_{p-q-1} = b_q/a_p \neq 0$ . If  $p - q \leq n$ , let  $C \in K[X], T \in K[[1/X]]$  be defined by  $S(C+T) = 1$ . It is clear that  $\deg(C) = p - q, P = CQ + R$ , with  $\deg(R) < q$  and

$$T = -R/Q = \sum_{j=0}^{\infty} t_j / X^{j+1}.$$

By Proposition 9.20, the matrix of coefficients of  $\text{Bez}(P, Q)$  in the Horner basis  $\text{Hor}_{p-1}(P, X), \dots, \text{Hor}_0(P, X)$  is  $\overline{\text{Han}}(\bar{s}_p)$ .

We consider now the matrix of coefficients of  $\text{Bez}(P, Q)$  in the basis

$$X^{p-q-1}Q(X), \dots, XQ(X), Q(X), \text{Hor}_{q-1}(Q, X), \dots, \text{Hor}_0(Q, X).$$

Since  $P = CQ + R$ ,  $\deg(R) < q$ ,

$$\text{Bez}(P, Q) = \frac{C(X) - C(Y)}{X - Y} Q(Y)Q(X) + \text{Bez}(Q, -R),$$

the matrix of coefficients of  $\text{Bez}(P, Q)$  in the basis

$$X^{p-q-1}Q(X), \dots, XQ(X), Q(X), \text{Hor}_{q-1}(Q, X), \dots, \text{Hor}_0(Q, X),$$



is the block Hankel matrix

$$\begin{bmatrix} \text{Han}(\bar{u}) & 0 \\ 0 & \text{Han}(\bar{t}_q) \end{bmatrix}.$$

**Proof of Lemma 9.24:** Take

$$P = X^{2n-1}, Q = s_{c-1}X^{2n-c-1} + \dots + s_{2n-2}.$$

Note that

$$Q/P = \sum_{j=c-1}^{\infty} s_j/X^{j+1} \in K[[1/X]]$$

with  $s_j = 0$  for  $j > 2n - 2$ . The Horner basis of  $P$  is  $X^{n-2}, \dots, X, 1$ . The change of basis matrix from the Horner basis of  $P$  to

$$X^{c-1}Q(X), \dots, XQ(X), Q(X), \text{Hor}_{q-1}(Q, X), \dots, \text{Hor}_0(Q, X)$$

is  $\text{To}(\bar{w})$ , with  $\bar{w} = s_{c-1}, s_c, \dots, s_{2n+c-3}$ . Thus, according to the preceding discussion,

$$\text{Han}(\bar{s}_{2n-1}) = \text{To}(\bar{w})^t \begin{bmatrix} \text{Han}(\bar{u}) & 0 \\ 0 & \text{Han}(\bar{t}_{2n-c-1}) \end{bmatrix} \text{To}(\bar{w}). \tag{9.11}$$

Lemma 9.24 follows easily from Equation (9.11), suppressing the last  $n - 1$  lines and columns in each matrix.  $\square$

The following result gives a method to compute the signature of a Hankel form.

**Proposition 9.25.** *Let*

$$\begin{aligned} P &= a_p X^p + \dots + a_0 \\ Q &= b_q X^q + \dots + b_0 \end{aligned}$$

*be coprime,  $q = \deg(Q) < p = \deg(P)$ , such that*

$$Q/P = \sum_{j=0}^{\infty} s_j/X^{j+1} \in K[[1/X]].$$

*Let  $\text{han}(\bar{s}_{[0..n]}) = 1, \text{han}(\bar{s}_1), \dots, \text{han}(\bar{s}_n)$ .*

*a) Suppose  $p \leq n$ . Then*

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{PmV}(\text{han}(\bar{s}_{[0..n]})) = \text{PmV}(\text{sRes}(P, Q)) = \text{Ind}(Q/P).$$

*b) Suppose  $p > n$ . Denoting by  $j$  the biggest natural number  $\leq p - n$  such that the subresultant  $\text{sResP}_j$  is non-defective and by  $\text{sRes}(P, Q)_{[p..j]}$  the sequence of  $\text{sRes}_i(P, Q)$ ,  $i = p, \dots, j$ , we have*

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{PmV}(\text{han}(\bar{s}_{[0..p-j]})) = \text{PmV}(\text{sRes}(P, Q)_{[p..j]}).$$

The proof of the proposition uses the following lemma

**Lemma 9.26.** *For all  $k \in \{1, \dots, p\}$ , we have:*

$$\text{sRes}_{p-k}(P, Q) = a_p^{2k+2-p+q} \text{han}(\bar{s}_k).$$

**Proof:** Let

$$\Delta = \begin{bmatrix} a_p & a_{p-1} & & & & \ddots & & \ddots & & a_{p-2k+2} \\ 0 & a_p & a_{p-1} & & & \ddots & & \ddots & & a_{p-2k+3} \\ \vdots & \ddots & \ddots & \ddots & & \ddots & & \ddots & & \vdots \\ \vdots & & & a_p & a_{p-1} & & & & & a_{p-k} \\ \vdots & & & 0 & b_{p-1} & & & & & b_{p-k} \\ \vdots & & & & & \ddots & & & & \vdots \\ 0 & b_{p-1} & & & & \ddots & & & & b_{p-2k+2} \\ b_{p-1} & & & & \ddots & & & \ddots & & b_{p-2k+1} \end{bmatrix},$$

(the first coefficients of  $Q$  may be zero).

We have  $\det(\Delta) = a_p^{p-q-1} \text{sRes}_{p-k}(P, Q)$ . From the relations (9.6), we deduce  $a_p^{p-q-1} \text{sRes}_{p-k}(P, Q) = \det(D) \det(D')$  with

$$D = \begin{bmatrix} 1 & 0 & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & 0 \\ 0 & 1 & 0 & & & & & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & & & & & & \vdots \\ \vdots & & & 1 & 0 & \cdots & \cdots & \cdots & \cdots & 0 \\ \vdots & & & 0 & s_0 & s_1 & \cdots & \cdots & \cdots & s_{k-1} \\ \vdots & & & & & & \ddots & & & \vdots \\ \vdots & & & & & & & & & \vdots \\ 0 & s_0 & s_1 & & & & \ddots & & & s_{2k-3} \\ s_0 & s_1 & \cdots & \cdots & \cdots & s_{k-1} & & \ddots & & s_{2k-2} \end{bmatrix},$$

$$D' = \begin{bmatrix} a_p & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & a_{p-2k+2} \\ 0 & a_p & \ddots & & & & & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & & & & & & \vdots \\ 0 & & 0 & a_p & \ddots & & & & & \vdots \\ 0 & & & 0 & a_p & \ddots & & & & \vdots \\ \vdots & & & & & \ddots & & & & \vdots \\ \vdots & & & & & & \ddots & & & \vdots \\ \vdots & & & & & & & & & \vdots \\ 0 & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & 0 & a_p & \vdots \end{bmatrix}.$$

This implies that

$$a_p^{p-q-1} \text{sRes}_{p-k}(P, Q) = a_p^{2k+1} \text{han}(\bar{s}_k). \quad \square$$

**Proof of Proposition 9.25:** a) It follows from Lemma 9.26 and Theorem 4.31 that

$$\text{PmV}(\text{han}(\bar{s}_{[0..n]})) = \text{PmV}(\text{sRes}(P, Q)) = \text{Ind}(Q/P).$$

So it remains to prove that

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{Ind}(Q/P).$$

The proof is by induction on the number of elements  $m$  of the euclidean remainder sequence of  $P$  and  $Q$ .

If  $m = 2$ , then  $Q = b$  is a constant, and the equality

$$\text{Ind}(b/P) = \begin{cases} \text{sign}(a_p b) & \text{if } p \text{ is odd,} \\ 0 & \text{if } p \text{ is even,} \end{cases}$$

is part of the proof of Theorem 4.31. The equality

$$\text{Sign}(\overline{\text{Han}}(s_n)) = \begin{cases} \text{sign}(a_p b) & \text{if } p \text{ is odd,} \\ 0 & \text{if } p \text{ is even,} \end{cases}$$

follows from Lemma 9.22, since here  $p \leq n$ ,  $C = P/b$ ,  $c = p$ ,  $T = 0$ ,  $s_{p-1} = b/a_p$ . Thus the theorem is true when  $m = 2$ .

If  $m > 2$ , the theorem follows by induction from Lemma 9.4 and Lemma 9.22, since  $\deg(C) = p - q$ ,  $s_{p-q-1} = b_q a_p \neq 0$ ,

$$T = -R/Q = \sum_{j=0}^{\infty} t_j/X^{j+1},$$

and the signed remainder sequence of  $Q$  and  $-R$  has  $m - 1$  elements.

b) It follows from Lemma 9.26 that

$$\text{PmV}(\text{han}(\bar{s}_{[0..p-j]})) = \text{PmV}(\text{sRes}(P, Q)_{[p..j]}).$$

So it remains to prove that

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{PmV}(\text{sRes}(P, Q)_{[p..j]}).$$

The proof is by induction on the number of elements  $m$  of non-zero elements in  $\text{sRes}(P, Q)_{[p..j]}$ .

If  $m = 2$ , then  $n < p - q$ ,

$$\text{PmV}(\text{sRes}_p(P, Q), 0, \dots, 0, \text{sRes}_q(P, Q)) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p - q \text{ is odd,} \\ 0 & \text{if } p - q \text{ is even,} \end{cases}$$

according to the definitions of  $\text{sRes}_p(P, Q)$ ,  $\text{sRes}_q(P, Q)$  (see Notation 8.29), and  $D$  (see Notation 9.1). The equality

$$\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \begin{cases} \text{sign}(a_p b_q) & \text{if } p - q \text{ is odd,} \\ 0 & \text{if } p - q \text{ is even,} \end{cases}$$

is a particular case of Lemma 9.5 since  $s_{p-q-1} = b_q/a_p \neq 0$ . Thus the theorem is true when  $n = 2$ .

If  $m > 2$ , the theorem follows by induction from Lemma 9.4 and Lemma 9.22, since  $\text{sRes}(Q, -R)_{[q..j]}$  has  $m - 1$  non-zero elements by Proposition 8.35. □

*Remark 9.27.* a) Proposition 9.25 is a generalization of Theorem 4.57, taking for  $Q$  the remainder of  $P'Q$  divided by  $P$ .

b) Note than given any

$$\bar{s}_n = s_0, \dots, s_{m-1}, 0, \dots, 0$$

such that  $s_{m-1} \neq 0$ ,  $P = T^m$  and  $Q = s_0 T^{m-1} + \dots + s_{m-1}$  satisfy the hypotheses of Proposition 9.25. Thus Proposition 9.25 provides a general method for computing the signature of a Hankel form.

c) Note also that when the Hankel matrix  $\text{Han}(\bar{s}_n)$  is invertible,  $p = n$  and  $\text{Sign}(\overline{\text{Han}}(\bar{s}_n)) = \text{PmV}(\text{han}(\bar{s}_{[0..n]}))$ . The signature of a quadratic form associated to an invertible Hankel matrix is thus determined by the signs of the principal minors of its associated matrix. This generalizes Remark 4.59. □

We are now ready to describe an algorithm computing the signature of a Hankel quadratic form. The complexity of this algorithm will turn out to be better than the complexity of the algorithm computing the signature of a general quadratic form (Algorithm 8.18 (Signature through Descartes)), because the special structure of a Hankel matrix is taken into account in the computation.

*Algorithm 9.6. [Signature of Hankel Form]*

- **Structure:** an ordered integral domain  $D$ .
- **Input:**  $2n - 1$  numbers  $\bar{s}_n = s_0, \dots, s_{2n-2}$  in  $D$ .
- **Output:** the signature of the Hankel quadratic form  $\overline{\text{Han}}(\bar{s}_n)$ .
- **Complexity:**  $O(n^2)$ .
- **Procedure:**
  - If  $s_i = 0$  for every  $i = 0, \dots, 2n - 2$ , output 0.
  - If  $s_i = 0$  for every  $i = 0, \dots, c - 2$ ,  $s_{c-1} \neq 0$ ,  $c \geq n$ , output

$$\begin{cases} \text{sign}(s_{c-1}) & \text{if } c \text{ is odd,} \\ 0 & \text{if } c \text{ is even.} \end{cases}$$

- Otherwise, let  $m$ ,  $0 < m \leq 2n - 2$ , be such that  $s_{m-1} \neq 0$ ,  $s_i = 0, i \geq m$ .
- If  $m \leq n$ , output

$$\begin{cases} \text{sign}(s_{m-1}) & \text{if } m \text{ is odd,} \\ 0 & \text{if } m \text{ is even.} \end{cases}$$

- If  $m > n$ , take  $P := T^m$ ,  $Q := s_0 T^{m-1} + \dots + s_{m-1}$  and apply a modified version of Algorithm 8.21 (Signed Subresultant) to  $P$  and  $Q$  stopping at the first non-defective  $\text{sResP}_j(P, Q)$  such that  $j \leq m - n$ . Compute  $\text{PmV}(\text{sRes}(P, Q)_{[m..j]})$ .

**Proof of correctness:** Use Lemma 9.5, Proposition 9.25 together with Remark 9.27 a).  $\square$

**Complexity analysis:** The complexity of this algorithms is  $O(n^2)$ , by the complexity analysis of Algorithm 8.21 (Signed Subresultant).

When  $s_0, \dots, s_{2n-2}$  are in  $\mathbb{Z}$ , of bitsizes bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O((\tau + \nu)n)$  where  $\nu$  is the bitsize of  $n$  according to Proposition 8.44.  $\square$

*Remark 9.28.* It is possible to compute the signature of a Hankel matrix with complexity  $\tilde{O}(n\tau)$  and binary complexity  $\tilde{O}(n^2\tau)$  using Remark 9.2.  $\square$

### 9.3 Number of Complex Roots with Negative Real Part

So far Cauchy index was used only for the computation of Tarski-queries. We describe an important application of Cauchy index to the determination of the number of complex roots with negative real part of a polynomial with real coefficients.

Let  $P(X) = a_p X^p + \dots + a_0 \in \mathbb{R}[X]$ ,  $a_p \neq 0$ , where  $\mathbb{R}$  is real closed, and  $\mathbb{C} = \mathbb{R}[i]$  is algebraically closed. Define  $F(X), G(X)$  by

$$P(X) = F(X^2) + XG(X^2). \tag{9.12}$$

Note that if  $p = 2m$  is even

$$\begin{aligned} F &= a_{2m} X^m + a_{2m-2} X^{m-1} + \dots, \\ G &= a_{2m-1} X^{m-1} + a_{2m-3} X^{m-2} + \dots, \end{aligned}$$

and if  $p = 2m + 1$  is odd

$$\begin{aligned} F &= a_{2m} X^m + a_{2m-2} X^{m-1} + \dots, \\ G &= a_{2m+1} X^m + a_{2m-1} X^{m-1} + \dots. \end{aligned}$$

We are going to prove the following result.

**Theorem 9.29.** *Let  $n(P)$  be the difference between the number of roots of  $P$  with positive real parts and the number of roots of  $P$  with negative real parts.*

$$n(P) = \begin{cases} -\text{Ind}(G/F) + \text{Ind}(XG/F) & \text{if } p \text{ is even,} \\ -\text{Ind}(F/XG) + \text{Ind}(F/G) & \text{if } p \text{ is odd.} \end{cases}$$

This result has useful consequences in control theory. When considering a linear differential equation with coefficients depending on parameters  $a_i$ ,

$$a_p y^{(p)}(t) + a_{p-1} y^{(p-1)}(t) + \dots + a_0 y(t) = 0, a_p \neq 0, \tag{9.13}$$

it is important to determine for which values of the parameters all the roots of the characteristic equation

$$P = a_p X^p + a_{p-1} X^{p-1} + \cdots + a_0 = 0, a_p \neq 0, \quad (9.14)$$

have negative real parts. Indeed if  $x_1, \dots, x_r$  are the complex roots of  $P$  with respective multiplicities  $\mu_1, \dots, \mu_r$ , the functions

$$e^{x_i t}, \dots, t^{\mu_i-1} e^{x_i t}, i = 1, \dots, r$$

form a basis of solutions of Equation (9.13) and when all the  $x_i$  have negative real parts, all the solutions of Equation (9.13) tend to 0 as  $t$  tends to  $+\infty$ , for every possible initial value. This is the reason why the set of polynomials of degree  $p$  which have all their complex roots with negative real part is called the **domain of stability** of degree  $p$ .

We shall prove the following result, as a corollary of Theorem 9.29.

**Theorem 9.30. [Liénard/Chipart]** *The polynomial*

$$P = a_p X^p + \cdots + a_0, a_p > 0,$$

*belongs to the domain of stability of degree  $p$  if and only if  $a_i > 0$ ,  $i = 0, \dots, p$ , and*

$$\begin{cases} \text{sRes}_m(F, G) > 0, \dots, \text{sRes}_0(F, G) > 0 & \text{if } p = 2m \text{ is even,} \\ \text{sRes}_{m+1}(XG, F) > 0, \dots, \text{sRes}_0(XG, F) > 0 & \text{if } p = 2m + 1 \text{ is odd.} \end{cases}$$

As a consequence, we can decide whether or not  $P$  belongs to the domain of stability by testing the signs of some polynomial conditions in the  $a_i$ , without having to approximate the roots of  $P$ .

**Exercise 9.3.** Determine the conditions on  $a, b, c, d$  characterizing the polynomials  $P = aX^3 + bX^2 + cX + d$ , belonging to the domain of stability of degree 3.

The end of the section is devoted to the proof of Theorem 9.29 and Theorem 9.30.

Define  $A(X), B(X)$ ,  $\deg(A) = p$ ,  $\deg(B) < p$ , as the real and imaginary parts of  $(-i)^p P(iX)$ . In other words,

$$\begin{aligned} A &= a_p X^p - a_{p-2} X^{p-2} + \cdots, \\ B &= -a_{p-1} X^{p-1} + a_{p-3} X^{p-3} + \cdots, \end{aligned}$$

so that when  $p$  is even  $A$  is even and  $B$  is odd (resp. when  $p$  is odd  $A$  is odd and  $B$  is even). Note that, using the definition of  $F$  and  $G$  (see (9.12)),

– if  $p = 2m$ ,

$$A = (-1)^m F(-X^2), B = (-1)^m XG(-X^2). \quad (9.15)$$

– if  $p = 2m + 1$ ,

$$A = (-1)^m XG(-X^2), B = -(-1)^m F(-X^2). \tag{9.16}$$

We are going to first prove the following result.

**Theorem 9.31. [Cauchy]** *Let  $n(P)$  be the difference between the number of roots of  $P$  with positive real part and the number of roots of  $P$  with negative real part. Then,*

$$n(P) = \text{Ind}(B/A).$$

A preliminary result on Cauchy index is useful.

**Lemma 9.32.** *Denote by  $t \mapsto (A_t, B_t)$  a semi-algebraic and continuous map from  $[0, 1]$  to the set of pairs of polynomials  $(A, B)$  of  $\mathbb{R}[X]$  with  $A$  monic of degree  $p$ ,  $\deg(B) < p$  (identifying pairs of polynomials with their coefficients). Suppose that  $A_0$  has a root  $x$  of multiplicity  $\mu$  in  $(a, b)$  and no other root in  $[a, b]$ , and  $B_0$  has no root in  $[a, b]$ . Then, for  $t$  small enough,*

$$\text{Ind}((B_0/A_0; a, b) = \text{Ind}((B_t/A_t; a, b).$$

**Proof:** Using Theorem 5.12 (Continuity of roots), there are two cases to consider:

– If  $\mu$  is odd, the number  $n$  of roots of  $A_t$  in  $[a, b]$  with odd multiplicity is odd, and thus the sign of  $A_t$  changes  $n$  times while the sign of  $B_t$  is fixed, and hence for  $t$  small enough,

$$\text{Ind}(B_0/A_0; a, b) = \text{Ind}(B_t/A_t; a, b) = \text{sign}(A_0^{(\mu)}(0) B_0(x)).$$

– If  $\mu$  is even, the number of roots of  $A_t$  in  $[a, b]$  with odd multiplicity is even, and thus for  $t$  small enough,

$$\text{Ind}(B_0/A_0; a, b) = \text{Ind}(B_t/A_t; a, b) = 0.$$

□

**Proof of Theorem 9.31:** We can suppose without loss of generality that  $P(0) \neq 0$ .

If  $A$  and  $B$  have a common root  $a + ib$ ,  $a \in \mathbb{R}, b \in \mathbb{R}$ ,  $b - ia$  is a root of  $P$ .

– If  $b = 0$ ,  $ia$  and  $-ia$  are roots of  $P$ , and  $P = (X^2 + a^2)Q$ . Denoting

$$(-i)^{p-2}Q(iX) = C(X) + iD(X), \quad C \in \mathbb{R}[X], D \in \mathbb{R}[X],$$

we have

$$\begin{aligned} A &= (X^2 - a^2)C, \\ B &= (X^2 - a^2)D. \end{aligned}$$

– If  $b \neq 0$ ,  $b + ia$ ,  $b - ia$ ,  $-b + ia$ ,  $-b - ia$ , are roots of  $P$  and

$$P = (X^4 + 2(a^2 - b^2)X^2 + (a^2 + b^2)^2)Q.$$

– Denoting

$$(-i)^{p-4}Q(iX) = C(X) + iD(X), \quad C \in \mathbb{R}[X], D \in \mathbb{R}[X],$$

we have

$$\begin{aligned} A &= (X^4 - 2(a^2 - b^2)X^2 + (a^2 + b^2)^2)C, \\ B &= (X^4 - 2(a^2 - b^2)X^2 + (a^2 + b^2)^2)D. \end{aligned}$$

In both cases  $n(P) = n(Q)$ ,  $\text{Ind}(B/A) = \text{Ind}(D/C)$ .

So we can suppose without loss of generality that  $P$  has no two roots on the imaginary axis and no two roots with opposite real part, and  $A$  and  $B$  are coprime.

Let  $x_1 = a_1 + i b_1, \dots, x_r = a_r + i b_r$ , be the roots of  $P$  with multiplicities  $\mu_1, \dots, \mu_r$ ,  $c$  be a positive number smaller than the difference between two distinct absolute values of  $a_i$ ,  $M$  a positive number bigger than twice the absolute value of the  $b_i$ . Consider for  $t \in [0, 1]$ , and  $i = 1, \dots, r$ ,

$$x_{i,t} = (1 - t)x_i + t(a_i + c/Mb_i),$$

and the polynomial

$$P_t(X) = (X - x_{1,t})^{\mu_1} \dots (X - x_{r,t})^{\mu_r}.$$

Note that  $P_0 = P$ ,  $P_1$  has only real roots, and for every  $t \in [0, 1]$  no two roots with opposite real parts, and hence for every  $t \in [0, 1]$ , defining

$$(-i)^p P_t(iX) = A_t(X) + iB_t(X), \quad A_t \in \mathbb{R}[X], B_t \in \mathbb{R}[X],$$

the polynomial  $A_t$  and  $B_t$  are coprime. Thus  $\text{Res}(A_t, B_t) \neq 0$  and by Proposition 9.10, denoting by  $M(A_t, B_t)$  the matrix of coefficients of  $\text{Bez}(A_t, B_t)$  in the canonical basis  $X^{p-1}, \dots, 1$ ,  $\det(M(A_t, B_t)) \neq 0$ . Thus the rank of  $M(A_t, B_t)$  is constantly  $p$  as  $t$  varies in  $[0, 1]$ . Hence by Proposition 9.6 the signature of  $M(A_t, B_t)$  is constant as  $t$  varies in  $[0, 1]$ . We have proved that, for every  $t \in [0, 1]$ ,  $\text{Ind}(B_t/A_t) = \text{Ind}(B/A)$ .

So, we need only to prove the claim for a polynomial  $P$  with all roots real and no opposite roots. This is done by induction on the degree of  $P$ .

The claim is obvious for a polynomial of degree 1 since if  $P = X - a$ ,  $A = X$ , and  $B = a$ ,  $\text{Ind}(a/X)$  is equal to 1 when  $a > 0$  and  $-1$  when  $a < 0$ .

Suppose that the claim is true for every polynomial of degree  $< p$  and consider  $P$  of degree  $p$ . Let  $a$  be the root of  $P$  with minimum absolute value among the roots of  $P$  and  $P = (X - a)Q$ .

If  $a > 0$ , we are going to prove, denoting

$$(-i)^{p-1}Q(iX) = C(X) + iD(X), \quad C \in \mathbb{R}[X], D \in \mathbb{R}[X],$$



that

$$\text{Ind}(B/A) = \text{Ind}(D/C) + 1. \tag{9.17}$$

We define  $P_t = (X - t)Q$ ,  $t \in (0, a]$  and denote

$$(-i)^p P_t(iX) = A_t(X) + iB_t(X), A_t \in \mathbb{R}[X], B_t \in \mathbb{R}[X].$$

Note that  $P_a = P$ , and for every  $t \in (0, a]$ ,  $P_t$  has only real roots, no opposite roots, and  $A_t$  and  $B_t$  are coprime. Thus  $\text{Res}(A_t, B_t) \neq 0$  and by Proposition 9.10, denoting by  $M(A_t, B_t)$  be the matrix of coefficients of  $\text{Bez}(A_t, B_t)$  in the canonical basis  $X^{p-1}, \dots, 1$ ,  $\det(M(A_t, B_t)) \neq 0$ . Thus the rank of  $M(A_t, B_t)$  is constantly  $p$  as  $t$  varies in  $(0, a]$ . Thus by Proposition 9.6 the signature of  $M(A_t, B_t)$  is constant as  $t$  varies in  $(0, a]$ . We have proved that, for every  $t \in (0, a]$ ,

$$\text{Ind}(B_t/A_t) = \text{Ind}(B/A). \tag{9.18}$$

We now prove that

$$\text{Ind}(B_t/A_t) = \text{Ind}(D/C) + 1, \tag{9.19}$$

if  $t$  is small enough. Note that

$$\begin{aligned} A_t(X) + iB_t(X) &= (-i)^p P_t(iX) \\ &= (X + it)(-i)^{p-1} Q(iX) \\ &= (X + it)(C(X) + iD(X)), \\ A_t(X) &= XC(X) - tD(X), \\ B_t(X) &= XD(X) + tC(X). \end{aligned}$$

For  $t$  small enough,  $A_t$  is close to  $XC(X)$  and  $B_t$  close to  $XD(X)$ .

- If  $p$  is odd,  $C(0) \neq 0$ ,  $D(0) = 0$  since  $D$  is odd and  $C$  and  $D$  have no common root. For  $t$  small enough, using Theorem 5.12 (Continuity of roots),  $A_t$  has a simple root  $y$  close to 0. The sign of  $B_t(y)$  is the sign of  $tC(0)$ . Hence for  $[a, b]$  small enough containing 0, and  $t$  sufficiently small,

$$\text{Ind}(D/C; a, b) = 0, \text{Ind}(B_t/A_t; a, b) = 1.$$

- If  $p$  is even,  $C(0) = 0$ ,  $D(0) \neq 0$  since  $C$  is odd and  $C$  and  $D$  have no common root.
  - If  $C'(0)D(0) > 0$ , there is a jump from  $-\infty$  to  $+\infty$  in  $D/C$  at 0, and  $A_t(0)$  has two roots close to 0, one positive and one negative. Hence for  $[a, b]$  small enough containing 0, and  $t$  sufficiently small,

$$\text{Ind}(D/C; a, b) = 1, \text{Ind}(B_t/A_t; a, b) = 2.$$

- If  $C'(0)D(0) < 0$ , there is a jump from  $+\infty$  to  $-\infty$  in  $D/C$  at 0 and  $A_t(0)$  has no root close to 0. Hence for  $[a, b]$  small enough containing 0, and  $t$  sufficiently small,

$$\text{Ind}(D/C; a, b) = -1, \text{Ind}(B_t/A_t; a, b) = 0.$$

Using Lemma 9.32 at the neighborhood of non-zero roots of  $C$ , Equation (9.19) follows. Equation (9.17) follows from Equation (9.18) and Equation (9.19).

If  $a < 0$ , a similar analysis, left to the reader, proves that

$$\text{Ind}(B/A) = \text{Ind}(D/C) - 1. \quad \square$$

**Proof of Theorem 9.29:**

– If  $p = 2m$ , let

$$\varepsilon = \begin{cases} \text{sign}_{x < 0, x \rightarrow 0}(G(x)/F(x)) & \text{if } \lim_{x < 0, x \rightarrow 0} |G(x)/F(x)| = \infty, \\ 0 & \text{otherwise.} \end{cases}$$

Then, since by (9.15)  $A = (-1)^m F(-X^2)$ ,  $B = (-1)^m XG(-X^2)$ ,

$$\begin{aligned} \text{Ind}(B/A) &= \text{Ind}(XG(-X^2)/F(-X^2)) \\ &= \text{Ind}(XG(-X^2)/F(-X^2); -\infty, 0) \\ &\quad + \text{Ind}(XG(-X^2)/F(-X^2); 0, +\infty) + \varepsilon \\ &= 2 \text{Ind}(XG(-X^2)/F(-X^2); -\infty, 0) + \varepsilon \\ &= -2 \text{Ind}(G(-X^2)/F(-X^2); -\infty, 0) + \varepsilon \\ &= -2 \text{Ind}(G(X)/F(X); -\infty, 0) - \varepsilon \\ &= -(G(X)/F(X); -\infty, 0) + \text{Ind}(XG(X)/F(X); -\infty, 0) + \varepsilon \\ &= -\text{Ind}(G/F) + \text{Ind}(XG/F). \end{aligned}$$

– If  $p = 2m + 1$ , let

$$\varepsilon = \begin{cases} \text{sign}_{x < 0, x \rightarrow 0}(F(x)/G(x)) & \text{if } \lim_{x < 0, x \rightarrow 0} (F(x)/G(x)) \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Then, since by (9.16)  $A = (-1)^m XG(-X^2)$ ,  $B = -(-1)^m F(-X^2)$ ,

$$\begin{aligned} \text{Ind}(B/A) &= -\text{Ind}(F(-X^2)/XG(-X^2)) \\ &= -\text{Ind}(F(-X^2)/XG(-X^2); -\infty, 0) \\ &\quad - \text{Ind}(F(-X^2)/XG(-X^2); 0, +\infty) - \varepsilon \\ &= -2 \text{Ind}(F(-X^2)/XG(-X^2); -\infty, 0) - \varepsilon \\ &= -2 \text{Ind}(F(X)/XG(X); -\infty, 0) - \varepsilon \\ &= -\text{Ind}(F(X)/XG(X); -\infty, 0) \\ &\quad + \text{Ind}(F(X)/G(X); -\infty, 0) - \varepsilon \\ &= -\text{Ind}(F/XG) + \text{Ind}(F/G). \end{aligned}$$

This proves the theorem, using Theorem 9.31. □

**Proof of Theorem 9.30:** If

$$P = a_p X^p + \cdots + a_0, a_p > 0$$

belongs to the domain of stability of degree  $p$ , it is the product of  $a_p$ , polynomials  $X + u$  with  $u > 0 \in \mathbb{R}$ , and  $X^2 + sX + t$  with  $s > 0 \in \mathbb{R}, t > 0 \in \mathbb{R}$ , and hence all the  $a_i, i = 0, \dots, p$ , are strictly positive. Thus,  $F$  and  $G$  have no positive real root, and  $\text{sign}(F(0)G(0)) = \text{sign}(a_0 a_1) = 1$ . Hence,

– If  $p = 2m$  is even,

$$\begin{aligned} \text{Ind}(G/F) &= -\text{Ind}(XG/F), \\ -p &= -\text{Ind}(G/F) + \text{Ind}(XG/F) \\ m &= \text{Ind}(G/F), \end{aligned}$$

– If  $p = 2m + 1$  is odd,

$$\begin{aligned} \text{Ind}(F/XG) &= -\text{Ind}(F/G) + 1, \\ -p &= -\text{Ind}(F/XG) + \text{Ind}(F/G) \\ m + 1 &= \text{Ind}(F/XG). \end{aligned}$$

The proof of the theorem follows, using Theorem 4.31. □

## 9.4 Bibliographical Notes

The use of quadratic forms for studying the Cauchy index appears in [90]. Bezoutians were considered by Sylvester [153]. The signature of Hankel forms has been studied by Frobenius [61].

A survey of some results appearing in this chapter can be found in [99, 63]. However it seems that the link between the various approaches for computing the Cauchy index, using subresultants, is recent [70, 140].

The domain of stability has attracted much attention, notably by Routh [139], Hurwitz [93], and Liénart/Chipart [108].

## Real Roots

---

In Section 10.1 we describe classical bounds on the roots of polynomials. In Section 10.2 we study real roots of univariate polynomials by a method based on Descartes's law of sign and Bernstein polynomials. These roots are characterized by intervals with rational endpoints. The method presented works only for archimedean real closed fields. In the second part of the chapter we study exact methods working in general real closed fields. Section 10.3 is devoted to exact sign determination in a real closed field and Section 10.4 to characterization of roots in a real closed field.

Besides their aesthetic interest, the specific methods of Section 10.2 are important in practical computations. This is the reason why we describe them, though they are less general than the methods of the second part of the chapter.

### 10.1 Bounds on Roots

We have already used a bound on the roots of a univariate polynomial in Chapter 5 (see Proposition 5.9). The following classical bound will also be useful.

In this section, we consider a polynomial

$$P = a_p X^p + \cdots + a_q X^q, \quad p > q, \quad a_q a_p \neq 0,$$

with coefficients in an ordered field  $K$ , a real closed field  $R$  containing  $K$ , and  $C = R[i]$ .

**Notation 10.1. [Cauchy bound]** We denote

$$C(P) = \sum_{q \leq i \leq p} \left| \frac{a_i}{a_p} \right|,$$

$$c(P) = \left( \sum_{q \leq i \leq p} \left| \frac{a_i}{a_q} \right| \right)^{-1}.$$

□

**Lemma 10.2. [Cauchy]** *The absolute value of any root of  $P$  in  $\mathbb{R}$  is smaller than  $C(P)$ .*

**Proof:** Let  $x \in \mathbb{R}$  be a root of  $P = a_p X^p + \dots + a_q X^q$ ,  $p > q$ . Then

$$a_p x = - \sum_{q \leq i \leq p-1} a_i x^{i-p+1}.$$

If  $|x| \geq 1$  this gives

$$\begin{aligned} |a_p| |x| &\leq \sum_{q \leq i \leq p-1} |a_i| |x|^{i-p+1} \\ &\leq \sum_{q \leq i \leq p-1} |a_i|. \end{aligned}$$

Thus it is clear that  $|x| \leq C(P)$ .

If  $|x| \leq 1$ , we have  $|x| \leq 1 \leq C(P)$ , since  $C(P) \geq 1$ . □

Similarly, we have

**Lemma 10.3.** *The absolute value of any non-zero root of  $P$  in  $\mathbb{R}$  is bigger than  $c(P)$ .*

**Proof:** This follows from Lemma 10.2 by taking the reciprocal polynomial  $X^p P(1/X)$ . □

**Corollary 10.4.** *If  $P \in \mathbb{Z}[X]$  had degree at most  $p$ , coefficients of bit length at most  $\tau$ , and  $p$  has bitsize at most  $\nu$ , the absolute values of the roots of  $P$  in  $\mathbb{R}$  are bounded by  $2^{\tau+\nu}$ .*

**Proof:** Follows immediately from Lemma 10.2, since  $p+1 \leq 2^\nu$ . □

The following proposition will be convenient when the polynomials we consider depend on parameters.

**Notation 10.5. [Modified Cauchy bound]** We denote

$$\begin{aligned} C'(P) &= (p+1) \cdot \sum_{q \leq i \leq p} \frac{a_i^2}{a_p^2}, \\ c'(P) &= \left( (p+1) \cdot \sum_{q \leq i \leq p} \frac{a_i^2}{a_q^2} \right)^{-1}. \end{aligned}$$

□

**Lemma 10.6.** *The absolute value of any root of  $P$  in  $\mathbb{R}$  is smaller than  $C'(P)$ .*

**Proof:** Let  $x \in \mathbb{R}$  be a root of  $P = a_p X^p + \dots + a_q X^q$ ,  $p > q$ . Then

$$a_p x = - \sum_{q \leq i \leq p-1} a_i x^{i-p+1}.$$

If  $|x| \geq 1$ , this gives

$$\begin{aligned} (a_p x)^2 &= \left( \sum_{q \leq i \leq p-1} a_i x^{i-p-1} \right)^2 \\ &\leq (p+1) \left( \sum_{q \leq i \leq p-1} a_i^2 \right). \end{aligned}$$

Thus  $|x| \leq C'(P)$ . If  $|x| \leq 1$ , we have  $|x| \leq 1 \leq C'(P)$ , since  $C(P) \geq 1$ . □

**Lemma 10.7.** *The absolute value of any non-zero root of  $P$  in  $\mathbb{R}$  is bigger than  $c'(P)$ .*

**Proof:** This follows from Lemma 10.6 by taking the reciprocal polynomial  $X^p P(1/X)$ . □

Our next aim is to give a bound on the divisors of a polynomial with integer coefficients.

We are going to use the following notions. If

$$P = a_p X^p + \dots + a_0 \in \mathbb{C}[X], a_p \neq 0,$$

the **norm** of  $P$  is

$$\|P\| = \sqrt{|a_p|^2 + \dots + |a_0|^2}.$$

The **length** of  $P$  is

$$\text{Len}(P) = |a_p| + \dots + |a_0|.$$

If  $z_1, \dots, z_p$  are the roots of  $P$  in  $\mathbb{C}$  counted with multiplicity so that

$$P = a_p \prod_{i=1}^p (X - z_i), \tag{10.1}$$

the **measure** of  $P$  is

$$\text{Mea}(P) = |a_p| \prod_{i=1}^p \max(1, |z_i|).$$

These three quantities are related as follows

**Proposition 10.8.**

$$\text{Len}(P) \leq 2^p \text{Mea}(P).$$

**Proposition 10.9.**

$$\text{Mea}(P) \leq \|P\|.$$

**Proof of Proposition 10.8:** By Lemma 2.12,

$$a_{p-k} = (-1)^k \left( \sum_{1 \leq i_1 < \dots < i_k \leq p} z_{i_1} \dots z_{i_k} \right) a_p.$$

Thus  $|a_{p-k}| \leq \binom{p}{k} \text{Mea}(P)$ , and

$$\text{Len}(P) \leq \sum_{k=0}^p \binom{p}{k} \text{Mea}(P) = 2^p \text{Mea}(P). \quad \square$$

The proof of Proposition 10.9 relies on the following lemma.

**Lemma 10.10.** *If  $P \in \mathbb{C}[X]$  and  $\alpha \in \mathbb{C}$ , then*

$$\|(X - \alpha)P(X)\| = \|(\bar{\alpha}X - 1)P(X)\|.$$

**Proof:** We have

$$\begin{aligned} \|(X - \alpha)P(X)\|^2 &= \sum_{j=0}^{p+1} (a_{j-1} - \alpha a_j)(\bar{a}_{j-1} - \bar{\alpha} \bar{a}_j) \\ &= (1 + |\alpha|^2) \|P\|^2 - \sum_{j=0}^p (\alpha a_j \bar{a}_{j-1} + \bar{\alpha} \bar{a}_j a_{j-1}), \end{aligned}$$

where  $a_{-1} = a_{p+1} = 0$ , since

$$(a_{j-1} - \alpha a_j)(\bar{a}_{j-1} - \bar{\alpha} \bar{a}_j) = |a_{j-1}|^2 + |\alpha|^2 |a_j|^2 - (\alpha a_j \bar{a}_{j-1} + \bar{\alpha} \bar{a}_j a_{j-1}).$$

Similarly

$$\begin{aligned} \|(\bar{\alpha}X - 1)P(X)\|^2 &= \sum_{j=0}^{p+1} (\bar{\alpha} a_{j-1} - a_j)(\alpha \bar{a}_{j-1} - \bar{a}_j) \\ &= (1 + |\alpha|^2) \|P\|^2 - \sum_{j=0}^p (\alpha a_j \bar{a}_{j-1} + \bar{\alpha} \bar{a}_j a_{j-1}). \end{aligned} \quad \square$$

**Proof of Proposition 10.9:** Let  $z_1, \dots, z_k$  be the roots of  $P$  outside of the unit disk. Then, by definition,  $M(P) = |a_p| \prod_{i=1}^k |z_i|$ . We consider the polynomial

$$R = a_p \prod_{i=1}^k (\bar{z}_i X - 1) \prod_{i=k+1}^n (X - z_i) = b_p X^p + \dots + b_0.$$

Noting that  $|b_p| = \text{Mea}(P)$ , and applying Lemma 10.10  $k$  times, we obtain  $\|P\| = \|R\|$ . Since  $\|P\|^2 = \|R\|^2 \geq |b_p|^2 = \text{Mea}(P)^2$ , the claim is proved.  $\square$

**Proposition 10.11.** *If  $P \in \mathbb{Z}[X]$  and  $Q \in \mathbb{Z}[X]$ ,  $\deg(Q) = q$ ,  $Q$  divides  $P$ , then*

$$\begin{aligned} \text{Mea}(Q) &\leq \text{Mea}(P), \\ \text{Len}(Q) &\leq 2^q \|P\|. \end{aligned}$$

**Proof:** Since the leading coefficient of  $Q$  divides the leading coefficient of  $P$  and every root of  $Q$  is a root of  $P$ , it is clear that  $\text{Mea}(Q) \leq \text{Mea}(P)$ . The other part of the claim follows using Proposition 10.9 and Proposition 10.8.  $\square$

**Corollary 10.12.** *If  $P \in \mathbb{Z}[X]$  and  $Q \in \mathbb{Z}[X]$  divides  $P$  in  $\mathbb{Z}[X]$ , then the bitsize of any coefficient of  $Q$  is bounded by  $q + \tau + \nu$ , where  $\tau$  is a bound on the bitsizes of the coefficients of  $P$  and  $\nu$  is the bitsize of  $p + 1$ .*

**Proof:** Notice that  $\|P\| < (p + 1) 2^\tau$ ,  $2^{\tau-1} \leq \text{Len}(Q)$ , and apply Proposition 10.11.  $\square$

The preceding bound can be used to estimate the bitsizes of the coefficients of the separable part of a polynomial. The **separable part of  $P$**  is a separable polynomial with the same set of roots as  $P$  in  $\mathbb{C}$ . The separable part of  $P$  is unique up to a multiplicative constant.

**Lemma 10.13.** *The polynomial  $P/\text{gcd}(P, P')$  is the separable part of  $P$ .*

**Proof:** Decompose  $P$  as a product of linear factors over  $\mathbb{C}$ :

$$P = (X - z_1)^{\mu_1} \cdots (X - z_r)^{\mu_r},$$

with  $z_1, \dots, z_r$  distinct. Then since  $z_1, \dots, z_r$  are roots of  $P'$  of multiplicities  $\mu_1 - 1, \dots, \mu_r - 1$ ,

$$\begin{aligned} \text{gcd}(P, P') &= (X - z_1)^{\mu_1 - 1} \cdots (X - z_r)^{\mu_r - 1}, \\ P/\text{gcd}(P, P') &= (X - z_1) \cdots (X - z_r). \end{aligned}$$

is separable.  $\square$

More generally, it is convenient to consider the **gcd-free part of  $P$**  with respect to  $Q$ , which is the divisor  $D$  of  $P$  such that  $DQ = \text{lcm}(P, Q)$ . It is clear that  $D = P/\text{gcd}(P, Q)$ . The gcd-free part of  $P$  with respect to  $Q$  is unique up to a multiplicative constant.

The greatest common divisor of  $P$  and  $Q$  and the gcd-free part of  $P$  with respect to  $Q$  can be computed using Algorithm 8.22 (Extended Signed Subresultant).

**Proposition 10.14.** *If  $\text{deg}(\text{gcd}(P, Q)) = j$ , then  $\text{sRes}P_j(P, Q)$  is the greatest common divisor of  $P$  and  $Q$  and  $\text{sRes}V_{j-1}(P, Q)$  is the gcd-free part of  $P$  with respect to  $Q$ .*

**Proof:** Let  $j$  be the degree of  $\text{gcd}(P, Q)$ . According to Theorem 8.30,  $\text{sRes}P_j(P, Q)$  is a greatest common divisor of  $P$  and  $Q$ . Moreover, Theorem 8.30 implies that  $\text{sRes}P_{j-1}(P, Q) = 0$ . Since, by Proposition 8.38,

$$\text{sRes}U_{j-1}(P, Q) P + \text{sRes}V_{j-1}(P, Q) Q = \text{sRes}P_{j-1}(P, Q) = 0,$$



then  $\text{sRes}U_{j-1}(P, Q) P = -\text{sRes}V_{j-1}(P, Q) Q$  is a multiple of the least common multiple of  $P$  and  $Q$  and is of degree  $\geq p + q - j$ . On the other hand by Proposition 8.38 a),

$$\begin{aligned}\deg(\text{sRes}U_{j-1}(P, Q)) &\leq q - j, \\ \deg(\text{sRes}V_{j-1}(P, Q)) &\leq p - j.\end{aligned}$$

It follows immediately that  $\text{sRes}U_{j-1}(P, Q)$  is proportional to  $Q/\text{gcd}(P, Q)$  and  $\text{sRes}V_{j-1}(P, Q)$  is proportional to  $P/\text{gcd}(P, Q)$ .  $\square$

**Corollary 10.15.** *If  $\deg(\text{gcd}(P, P')) = j$ ,  $\text{sRes}P_j(P, P')$  is the greatest common divisor of  $P$  and  $P'$  and  $\text{sRes}V_{j-1}(P, P')$  is the separable part of  $P$ .*

According to the preceding results, we are going to compute the gcd and gcd-free part using Algorithm 8.22 (Extended Signed Subresultants). In the case of integer coefficients, it will be possible to improve slightly the algorithm, using the following definitions and results.

If  $P \in \mathbb{Z}[X]$ , denote by  $\text{cont}(P)$  the **content** of  $P$ , which is the greatest common divisor of the coefficients of  $P$ .

**Lemma 10.16.** *Let  $P_1 \in \mathbb{Z}[X]$ ,  $P_2 \in \mathbb{Z}[X]$ . If  $\text{cont}(P_1) = \text{cont}(P_2) = 1$ , then  $\text{cont}(P_1 P_2) = 1$ .*

**Proof:** Consider a prime number  $p$ . Reducing the coefficients of  $P_1$  and  $P_2$  modulo  $p$ , notice that if  $P_1$  and  $P_2$  are not zero modulo  $p$ ,  $P_1 P_2 = P$  is also not zero modulo  $p$ . Thus  $\text{cont}(P)$  is divisible by no prime number  $p$ , and hence is equal to 1.  $\square$

**Lemma 10.17.** *If  $P \in \mathbb{Z}[X]$ ,  $P = P_1 P_2$ ,  $P_1 \in \mathbb{Q}[X]$  and  $P_2 \in \mathbb{Q}[X]$  there exist  $\bar{P}_1$  and  $\bar{P}_2$  in  $\mathbb{Z}[X]$ , proportional to  $P_1$  and  $P_2$  resp., such that  $\bar{P}_1 \bar{P}_2 = P$ .*

**Proof:** We can easily find  $\bar{P}_1 \in \mathbb{Z}[X]$  proportional to  $P_1$  such that  $\text{cont}(\bar{P}_1) = 1$ . Let  $c$  be the least common multiple of the denominators of the coefficients of  $\bar{P}_2 = P/\bar{P}_1$ . Then  $c\bar{P}_2 \in \mathbb{Z}[X]$  and  $\text{cont}(c\bar{P}_2) = d$  is prime to  $c$ . Consider  $\bar{P}_1$  and  $c\bar{P}_2/d$ , which belong to  $\mathbb{Z}[X]$ . Both of these polynomials have content equal to 1 and hence  $\text{cont}(cP/d) = 1$ , by Lemma 10.16. Since  $c$  and  $d$  are coprime,  $c = 1$ ,  $\text{cont}(P) = d$ . Hence  $\bar{P}_2 \in \mathbb{Z}[X]$ .  $\square$

#### Algorithm 10.1. [Gcd and Gcd-free Part]

- **Structure:** an integral domain  $D$ .
- **Input:** two univariate polynomials  $P = a_p X^p + \dots + a_0$  and  $Q = b_q X^q + \dots + b_0$  with coefficients  $D$  and of respective degrees  $p$  and  $q$ ,  $p > q$ .
- **Output:** the greatest common divisor of  $P$  and  $Q$  and the gcd-free part of  $P$  with respect to  $Q$ .
- **Complexity:**  $O(p^2)$ .

- **Procedure:**
  - Run Algorithm 8.22 (Extended Signed Subresultants) with inputs  $P$  and  $Q$ . Let  $j = \deg(\gcd(P, Q))$ .
  - If  $D = \mathbb{Z}$ , output  $a_p \text{sRes}P_j/\text{sRes}_j, a_p \text{sRes}V_{j-1}/\text{lcof}(\text{sRes}V_{j-1})$ .
  - Otherwise, output  $\text{sRes}P_j, \text{sRes}V_{j-1}$ .

**Proof of correctness of Algorithm 10.1:** The correctness of the algorithm when  $D \neq \mathbb{Z}$  follows from the correctness of Algorithm 8.22 (Extended Signed Subresultants) and Corollary 10.15.

When  $D = \mathbb{Z}$ , Lemma 10.17 implies that there exists  $a$  in  $\mathbb{Z}$  with  $a$  dividing  $a_p$  such that  $a \text{sRes}P_j/\text{sRes}_j \in \mathbb{Z}[X]$  and there exists  $b$  in  $\mathbb{Z}$  with  $b$  dividing  $a_p$  such that  $b \text{sRes}V_{j-1}/\text{lcof}(\text{sRes}V_{j-1}) \in \mathbb{Z}[X]$ . Thus

$$a_p \text{sRes}P_j/\text{sRes}_j \in \mathbb{Z}[X],$$

$$a_p \text{sRes}V_{j-1}/\text{lcof}(\text{sRes}V_{j-1}) \in \mathbb{Z}[X]$$

□

**Complexity analysis:** The complexity is  $O(p^2)$ , using the complexity analysis of Algorithm 8.22 (Extended Signed Subresultants).

When  $P \in \mathbb{Z}[X]$ , with the bitsizes of its coefficients bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O(\tau p)$  according to Proposition 8.44. Moreover using Corollary 10.12 the bitsize of the output is  $j + \tau + \nu$  and  $p - j + \tau + \nu$  with  $\nu$  the bitsize of  $p + 1$ . Note that the bitsize produced by the subresultant algorithm would be  $(p + q - 2j)(\tau + \mu)$  with  $\mu$  the bitsize of  $p + q$ , so the normalization step at the end of the algorithm when  $D = \mathbb{Z}$  improves the bounds on the bitsizes of the result. □

*Remark 10.18.* Algorithm 10.1 (Gcd and Gcd-Free part) is based on the Algorithm 8.22 (Extended Signed Subresultants) which uses exact divisions and is valid only in an integral domain, and not in a general ring. In a ring, the algorithm for computing determinants indicated in Remark 8.19 can always be used for computing the signed subresultant coefficients, and hence the separable and the gcd-free part. The complexity is  $(pq)^{O(1)}$  arithmetic operations in the ring  $D$  of coefficients of  $P$  and  $Q$ , which is sufficient for the complexity estimates proved in later chapters. □

*Remark 10.19.* The computation of the gcd and gcd free-part can be performed with complexity  $\tilde{O}(d)$  using [145, 119], and with binary complexity  $\tilde{O}(d^2 \tau)$ , using [107]. □

Now we study the minimal distance between roots of a polynomial  $P$ .

If  $P = a_p \prod_{i=1}^p (X - z_i) \in \mathbb{C}[Z]$ , the **minimal distance between the roots of  $P$**  is

$$\text{sep}(P) = \min \{|z_i - z_j|, z_i \neq z_j\}.$$

We denote by  $\text{Disc}(P)$  the discriminant of  $P$  (see Notation 4.1).

**Proposition 10.20.**

$$\begin{aligned} \text{sep}(P) &\geq (p/\sqrt{3})^{-1} p^{-p/2} |\text{Disc}(P)|^{1/2} \text{Mea}(P)^{1-p} \\ &\geq (p/\sqrt{3})^{-1} p^{-p/2} |\text{Disc}(P)|^{1/2} \|P\|^{1-p}. \end{aligned}$$

**Proof:** Consider the Vandermonde matrix

$$V(z_1, \dots, z_p) = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ z_1 & z_2 & \cdots & z_p \\ \vdots & \vdots & & \vdots \\ z_1^{p-1} & z_2^{p-1} & \cdots & z_p^{p-1} \end{bmatrix}.$$

We know by Equation (4.4) that  $\text{Disc}(P) = a_p^{2p-2} \det(V(z_1, \dots, z_p))^2$ . We suppose without loss of generality that  $|z_1 - z_2| = \text{sep}(P)$  and  $|z_1| \geq |z_2|$ . Using Hadamard's inequality (Proposition 8.9) on

$$V' = \begin{bmatrix} 0 & 1 & \cdots & 1 \\ z_1 - z_2 & z_2 & \cdots & z_p \\ \vdots & \vdots & & \vdots \\ z_1^{p-1} - z_2^{p-1} & z_2^{p-1} & \cdots & z_p^{p-1} \end{bmatrix},$$

and noticing that  $\det(V(z_1, \dots, z_p)) = \det(V')$ , we get

$$|\text{Disc}(P)|^{1/2} \leq |a_p|^{p-1} \left( \sum_{j=0}^{p-1} |z_1^j - z_2^j|^2 \right)^{1/2} \prod_{i \neq 1} (1 + |z_i|^2 + \cdots + |z_i|^{2(p-1)})^{1/2}.$$

Now,

$$\begin{aligned} \prod_{i \neq 1} (1 + |z_i|^2 + \cdots + |z_i|^{2(p-1)})^{1/2} &\leq \prod_{i \neq 1} (p \max(1, |z_i|)^{2(p-1)})^{1/2} \\ &\leq p^{(p-1)/2} \left( \frac{\text{Mea}(P)}{|a_p| \max(1, |z_1|)} \right)^{p-1}. \end{aligned}$$

On the other hand since it is clear that  $\sum_{j=0}^{p-1} j^2 \leq p^3/3$ ,

$$\begin{aligned} |z_1^j - z_2^j| &\leq j |z_1 - z_2| \max(1, |z_1|)^{p-1}, \\ \sum_{j=0}^{p-1} |z_1^j - z_2^j|^2 &\leq \left( \sum_{j=0}^{p-1} j^2 \right) |z_1 - z_2|^2 \max(1, |z_1|)^{2p-2}, \\ &\leq (p^3/3) |z_1 - z_2|^2 \max(1, |z_1|)^{2p-2}, \\ \left( \sum_{j=0}^{p-1} |z_1^j - z_2^j|^2 \right)^{1/2} &\leq (p^{3/2}/\sqrt{3}) |z_1 - z_2| \max(1, |z_1|)^{p-1}. \end{aligned}$$

Finally

$$|\text{Disc}(P)|^{1/2} \leq \text{sep}(P) (p/\sqrt{3}) p^{p/2} \text{Mea}(P)^{p-1}. \quad \square$$

**Proposition 10.21.** *If  $P \in \mathbb{Z}[X]$ ,*

$$\text{sep}(P) \geq (p/\sqrt{3})^{-1} p^{-p/2} \text{Mea}(P)^{1-p} \geq (p/\sqrt{3})^{-1} p^{-p/2} \|P\|^{1-p}.$$

**Proof:** If  $P$  is separable,  $\text{Disc}(P)$  is a non-zero integer, by Proposition 4.2 and Remark 4.28. Hence  $|\text{Disc}(P)| \geq 1$  and the claim follows by Proposition 10.20.

If  $P$  is not separable, its separable part  $Q$  divides  $P$  and belongs to  $\mathbb{Z}[X]$ . Thus by Proposition 10.11,  $\text{Mea}(Q) \leq \text{Mea}(P)$ . The conclusion follows, using Proposition 10.20 for  $Q$  and  $|\text{Disc}(Q)| \geq 1$ . □

**Corollary 10.22.** *If  $P$  is of degree at most  $p$  with coefficients in  $\mathbb{Z}$  of bitsize bounded by  $\tau$*

$$\text{sep}(P) \geq (p/\sqrt{3})^{-1} p^{-p/2} (p+1)^{(1-p)/2} 2^{\tau(1-p)}.$$

**Proof:** It is clear that  $\|P\| \leq (p+1)^{1/2} 2^\tau$ . □

Proposition 10.20 can be generalized as follows.

**Proposition 10.23.** *Let  $Z = \text{Zer}(P, \mathbb{C})$  and  $G = (Z, E)$  a directed acyclic graph such that*

- a) *for all  $(z_i, z_j) \in E$ ,  $|z_i| \leq |z_j|$ ,*
- b) *the in-degree of  $G$  is at most 1, i.e. for every  $z_j \in Z$ , there is at most one element  $z_i$  of  $Z$  such that  $(z_i, z_j) \in E$ .*

*Then, denoting by  $m$  the number of elements of  $E$ ,*

$$\begin{aligned} \prod_{(z_i, z_j) \in E} |z_i - z_j| &\geq (p/\sqrt{3})^{-m} p^{-p/2} |\text{Disc}(P)|^{1/2} \text{Mea}(P)^{1-p} \\ &\geq (p/\sqrt{3})^{-m} p^{-p/2} |\text{Disc}(P)|^{1/2} \|P\|^{1-p}. \end{aligned}$$

**Proof:** We can suppose without loss of generality that  $(z_i, z_j) \in E$  implies  $j < i$ . Consider, as in the proof of Proposition 10.20, the Vandermonde matrix  $V(z_1, \dots, z_p)$ .

Denote by  $A$  the subset of  $\{1, \dots, p\}$  consisting of elements  $j$  such that there exists  $i$  (necessarily greater than  $j$ ) such that  $(z_i, z_j) \in E$ , and note that by condition (b), the number of elements of  $A$  is  $m$ . For  $j \in A$  in increasing order, replace the  $j$ -th column of  $V(z_1, \dots, z_p)$  by the difference between the  $j$ -th and  $i$ -th column and denote by  $V_G$  the corresponding matrix. Because of condition (b),  $\det(V(z_1, \dots, z_p)) = \det(V_G)$ .

Using Hadamard's inequality (Proposition 8.9) on  $V_G$  and using the fact that, by Equation (4.4),  $\text{Disc}(P) = a_p^{2p-2} \det(V(z_1, \dots, z_p))^2$

$$\begin{aligned} &|\text{Disc}(P)|^{1/2} \\ &\leq |a_p|^{p-1} \prod_{(z_i, z_j) \in E} \left( \sum_{k=0}^p |z_i^k - z_j^k|^2 \right)^{1/2} \prod_{j \notin A} (1 + |z_j|^2 + \dots + |z_j|^{2(p-1)})^{1/2}. \end{aligned}$$

Now,

$$\begin{aligned} \prod_{j \notin A} (1 + |z_j|^2 + \dots + |z_j|^{2(p-1)})^{1/2} &\leq \prod_{j \notin A} (p \max(1, |z_j|)^{2(p-1)})^{1/2} \\ &\leq p^{(p-m)/2} \left( \frac{\text{Mea}(P)}{|a_p| \prod_{j \in A} \max(1, |z_j|)} \right)^{p-1}. \end{aligned}$$

On the other hand since it is clear that  $\sum_{k=0}^{p-1} k^2 \leq p^3/3$ ,

$$\begin{aligned} |z_i^k - z_j^k| &\leq j |z_i - z_j| \max(1, |z_j|)^{p-1}, \\ \sum_{k=0}^{p-1} |z_i^k - z_j^k|^2 &\leq \left( \sum_{k=0}^{p-1} k^2 \right) |z_i - z_j|^2 \max(1, |z_j|)^{2p-2}, \\ &\leq (p^3/3) |z_i - z_j|^2 \max(1, |z_j|)^{2p-2}, \\ \left( \sum_{k=0}^{p-1} |z_i^k - z_j^k|^2 \right)^{1/2} &\leq (p^{3/2}/\sqrt{3}) |z_i - z_j| \max(1, |z_j|)^{p-1}, \\ \prod_{(z_i, z_j) \in E} \left( \sum_{k=0}^{p-1} |z_i^k - z_j^k|^2 \right)^{1/2} &\leq (p^{3/2}/\sqrt{3})^m B, \end{aligned}$$

with  $B = \prod_{(z_i, z_j) \in E} |z_i - z_j| \prod_{j \in A} \max(1, |z_j|)^{p-1}$ .

Finally

$$|\text{Disc}(P)|^{1/2} \leq \prod_{(z_i, z_j) \in E} |z_i - z_j| (p/\sqrt{3})^m p^{p/2} \text{Mea}(P)^{p-1}. \quad \square$$

**Corollary 10.24.** *Let  $P$  be of degree at most  $p$  with coefficients in  $\mathbb{Z}$  of bitsize bounded by  $\tau$ . Let  $Z = \text{Zer}(P, \mathbb{C})$  and  $G = (Z, E)$  a directed acyclic graph such that*

- a) for all  $(z_i, z_j) \in E$ ,  $|z_i| \leq |z_j|$ ,
- b) the in-degree of  $G$  is at most 1.

Then, denoting by  $m$  the number of elements of  $E$ ,

$$\prod_{(z_i, z_j) \in E} |z_i - z_j| \geq (p/\sqrt{3})^{-m} p^{-p/2} (p+1)^{(1-p)/2} 2^\tau(1-p).$$

**Proof:** It is clear that  $\|P\| \leq (p+1)^{1/2} 2^\tau$ . □

## 10.2 Isolating Real Roots

Throughout this section,  $\mathbb{R}$  is an archimedean real closed field. Let  $P$  be a polynomial of degree  $p$  in  $\mathbb{R}[X]$ . We are going to explain how to perform exact computations for determining several properties of the roots of  $P$  in  $\mathbb{R}$ : characterization of a root, sign of another polynomial at a root, and comparisons between roots of two polynomials.

The characterization of the roots of  $P$  in  $\mathbb{R}$  will be performed by finding intervals with rational end points. Our method will be based on Descartes’s law of signs (Theorem 2.33) and the properties of the Bernstein basis defined below.

**Notation 10.25. [Bernstein polynomials]** The **Bernstein polynomials** of degree  $p$  for  $\ell, r$  are

$$\text{Bern}_{p,i}(\ell, r) = \binom{p}{i} \frac{(X - \ell)^i (r - X)^{p-i}}{(r - \ell)^p},$$

for  $i = 0, \dots, p$ . □

*Remark 10.26.* Note that  $\text{Bern}_{p,i}(\ell, r) = (-1)^p \text{Bern}_{p,p-i}(r, \ell)$  and that

$$\begin{aligned} \text{Bern}_{p,i}(\ell, r) &= \frac{(X - \ell)}{r - \ell} \frac{p}{i} \text{Bern}_{p-1,i-1}(\ell, r) \\ &= \frac{(r - X)}{r - \ell} \frac{p}{p-i} \text{Bern}_{p-1,i}(\ell, r). \end{aligned}$$

□

In order to prove that the Bernstein polynomials form a basis of polynomials of degree  $\leq p$ , we are going to need three simple transformations of  $P$ .

– **Reciprocal polynomial in degree  $p$ :**

$$\text{Rec}_p(P(X)) = X^p P(1/X).$$

The non-zero roots of  $P$  are the inverses of the non-zero roots of  $\text{Rec}(P)$ .

– **Contraction by ratio  $\lambda$ :** for every non-zero  $\lambda$ ,

$$\text{Co}_\lambda(P(X)) = P(\lambda X).$$

The roots of  $\text{Co}_\lambda(P)$  are of the form  $x/\lambda$ , where  $x$  is a root of  $P$ .

– **Translation by  $c$ :** for every  $c$ ,

$$\text{T}_c(P(X)) = P(X - c).$$

The roots of  $\text{T}_c(P(X))$  are of the form  $x + c$  where  $x$  is a root of  $P$ .

These three transformations clearly define bijections from the set of polynomials of degree at most  $p$  into itself.

**Proposition 10.27.** *Let  $P = \sum_{i=0}^p b_i \text{Bern}_{p,i}(\ell, r) \in \mathbb{R}[X]$  be of degree  $\leq p$ . Let  $\text{T}_{-1}(\text{Rec}_p(\text{Co}_{r-\ell}(\text{T}_{-\ell}(P)))) = \sum_{i=0}^p c_i X^i$ . Then  $\binom{p}{i} b_i = c_{p-i}$ .*

**Proof:** Performing the contraction of ratio  $r - \ell$  after translating by  $-\ell$  transforms  $\binom{p}{i} (X - \ell)^i (r - X)^{p-i} / (r - \ell)^p$  into  $\binom{p}{i} X^i (1 - X)^{p-i}$ . Translating by  $-1$  after taking the reciprocal polynomial in degree  $p$  transforms  $\binom{p}{i} X^i (1 - X)^{p-i}$  into  $\binom{p}{i} X^{p-i}$ . □

*Remark 10.28.* Proposition 10.27 immediately provides an algorithm for converting a polynomial from the monomial basis to the Bernstein basis for  $\ell, r$ .  $\square$

**Corollary 10.29.** *The Bernstein polynomials for  $c, d$  form a basis of the vector space of polynomials of degree  $\leq p$ .*

**Corollary 10.30.** *Let  $P \in \mathbb{Z}[X]$  be of degree  $\leq p$ . If the bitsizes of the coefficients of  $P$  are bounded by  $\tau$  in the monomial basis  $1, X, \dots, X^{p-1}$  and the bitsizes of the rational numbers  $r$  and  $\ell$  are bounded by  $\tau'$ , then there exists  $\lambda(P, \ell, r) \in \mathbb{Z}$  such that the bitsizes of the coefficients of  $\lambda(P, \ell, r)P$  in the Bernstein basis for  $(r, \ell)$  are integers of bitsize bounded by  $O(\tau + p\tau' + p \log_2(p))$ .*

**Proof:** Let  $\ell = a/b, r - \ell = a'/b'$ , with  $a, b, a', b'$  in  $\mathbb{Z}$ . Consider

$$p! T_{-1}(\text{Rec}_p(b'^p \text{Co}_{a'/b'}(b^p T_{-c}(P)))) = \sum_{i=0}^p d_i X^i.$$

It is clear that  $d_i$  is an integer multiple of  $p!$ . Thus the quotient  $\bar{b}_i$  of  $d_i$  by  $\binom{p}{i}$  is an integer, and we obtained  $\lambda(P)P$  with integer coefficients in the Bernstein basis of  $(c, d)$ . The claim follows immediately from Proposition 10.27 and the complexity analysis of Algorithm 8.10 (Special translation), noting that the bitsize of  $p!$  is  $O(p \log_2(p))$  by Stirling's formula.  $\square$

*Remark 10.31.* The list  $b = b_0, \dots, b_p$  of coefficients of  $P$  in the Bernstein basis of  $r, \ell$  gives the value of  $P$  at  $\ell$  (resp.  $r$ ), which is equal to  $b_0$  (resp.  $b_p$ ). Moreover, the sign of  $P$  at the right of  $\ell$  (resp. left of  $r$ ) is given by the first non-zero element (resp. last non-zero element) of the list  $b$ .  $\square$

We denote as usual by  $\text{Var}(b)$  the number of sign variations in a list  $b$ .

Note that, if  $\text{Var}(b) = 0$ , where  $b = b_0, \dots, b_p$  is the list of coefficients of  $P$  in the Bernstein basis of  $\ell, r$ , the sign of  $P$  on  $(c, d)$  is the sign of any non zero element of  $b$ , since the Bernstein polynomials for  $\ell, r$  are positive on  $(\ell, r)$ , thus  $P$  has no root in  $(\ell, r)$ . More generally, we have:

**Proposition 10.32.** *Let  $P$  be of degree  $p$ . We denote by  $b = b_0, \dots, b_p$  the coefficients of  $P$  in the Bernstein basis of  $\ell, r$ . Let  $\text{num}(P; (\ell, r))$  be the number of roots of  $P$  in  $(\ell, r)$  counted with multiplicities. Then*

- $\text{Var}(b) \geq \text{num}(P; (\ell, r))$ ,
- $\text{Var}(b) - \text{num}(P; (\ell, r))$  is even.

**Proof:** The claim follows immediately from Descartes's law of signs (Theorem 2.33), using Proposition 10.27. Indeed, the image of  $(\ell, r)$  under translation by  $-\ell$  followed by contraction of ratio  $r - \ell$  is  $(0, 1)$ . The image of  $(0, 1)$  under the inversion  $z \mapsto 1/z$  is  $(1, +\infty)$ . Finally, translating by  $-1$  gives  $(0, +\infty)$ .  $\square$

The coefficients  $b = b_0, \dots, b_p$  of  $P$  in the Bernstein basis of  $\ell, r$  give a rough idea of the shape of the polynomial  $P$  on the interval  $[\ell, r]$ . The **control line of  $P$  on  $[\ell, r]$**  is the union of the segments  $[M_i, M_{i+1}]$  for  $i = 0, \dots, p - 1$ , with

$$M_i = \left( \frac{i r + (p - i) \ell}{p}, b_i \right).$$

It is clear from the definitions that the graph of  $P$  goes through  $M_0$  and  $M_p$  and that the line  $M_0, M_1$  (resp  $M_{p-1}, M_p$ ) is tangent to the graph of  $P$  at  $M_0$  (resp.  $M_p$ ).

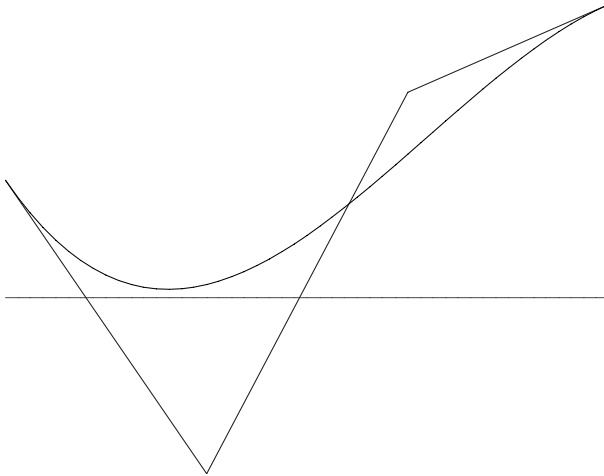
*Example 10.33.* We take  $p = 3$ , and consider the polynomial

$$P = -33 X^3 + 69 X^2 - 30 X + 4$$

with coefficients  $(4, -6, 7, 10)$  in the Bernstein basis for  $0, 1$

$$(1 - X)^3, 3 X (1 - X)^2, 3 X^2 (1 - X), X^3.$$

In Figure 10.1 we depict the graph of  $P$  on  $[0, 1]$ , the control line, and the  $X$ -axis.



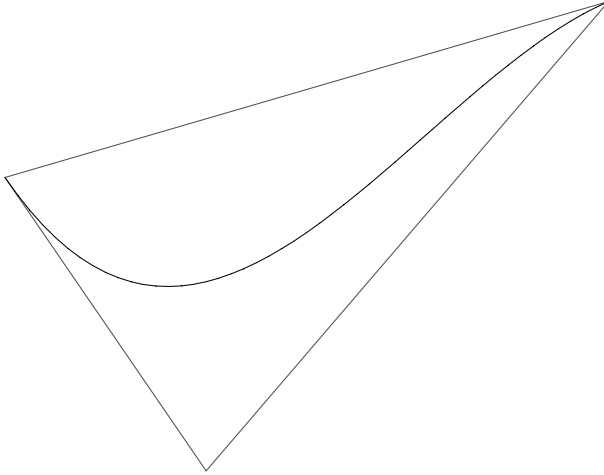
**Fig. 10.1.** Graph of  $P$  and control line of  $P$  on  $[0, 1]$

□

The **control polygon of  $P$  on  $[\ell, r]$**  is the convex hull of the points  $M_i$  for  $i = 0, \dots, p$ .

*Example 10.34.* Continuing Example 10.33, we draw the graph of  $P$  on  $[0, 1]$  and the control polygon (see Figure 10.2).





**Fig. 10.2.** Graph of  $P$  on  $[0, 1]$  and the control polygon

□

**Proposition 10.35.** *The graph of  $P$  on  $[\ell, r]$  is contained in the control polygon of  $P$  on  $[\ell, r]$ .*

**Proof:** In order to prove the proposition, it is enough to prove that any line  $L$  above (resp. under) all the points in the control polygon of  $P$  on  $[\ell, r]$  is above (resp. under) the graph of  $P$  on  $[\ell, r]$ . If  $L$  is defined by  $Y = aX + b$ , let us express the polynomial  $aX + b$  in the Bernstein basis. Since

$$1 = \left( \frac{X - \ell}{r - \ell} + \frac{r - X}{r - \ell} \right)^p,$$

the binomial formula gives

$$\begin{aligned} 1 &= \sum_{i=0}^p \binom{p}{i} \left( \frac{X - \ell}{r - \ell} \right)^i \left( \frac{r - X}{r - \ell} \right)^{p-i} \\ &= \sum_{i=0}^p \text{Bern}_{p,i}(\ell, r). \end{aligned}$$

Since

$$X = \left( r \left( \frac{X - \ell}{r - \ell} \right) + \ell \left( \frac{r - X}{r - \ell} \right) \right) \left( \frac{X - \ell}{r - \ell} + \frac{r - X}{r - \ell} \right)^{p-1},$$

the binomial formula together with Remark 10.31 gives

$$\begin{aligned} X &= \sum_{i=0}^{p-1} \left( r \left( \frac{X - \ell}{r - \ell} \right) + \ell \left( \frac{r - X}{r - \ell} \right) \right) \text{Bern}_{p-1,i}(\ell, r), \\ &= \sum_{i=0}^p \left( \frac{i r + (p - i) \ell}{p} \right) \text{Bern}_{p,i}(\ell, r). \end{aligned}$$

Thus,

$$aX + b = \sum_{i=0}^p \left( a \left( \frac{ir + (p-i)\ell}{p} \right) + b \right) \text{Bern}_{p,i}(\ell, r).$$

It follows immediately that if  $L$  is above every  $M_i$ , i.e. if

$$a \left( \frac{ir + (p-i)\ell}{p} \right) + b \geq b_i$$

for every  $i$ , then  $L$  is above the graph of  $P$  on  $[\ell, r]$ , since  $P = \sum_{i=0}^p b_i \text{Bern}_{p,i}(\ell, r)$  and the Bernstein basis of  $\ell, r$  is non-negative on  $[\ell, r]$ . A similar argument holds for  $L$  under every  $M_i$ . □

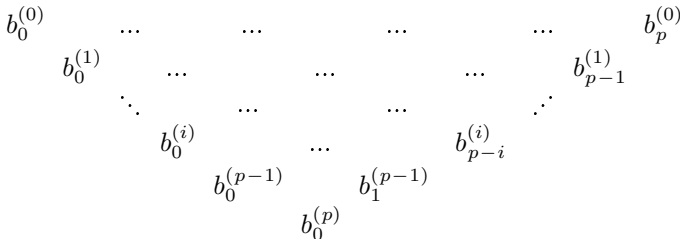
The following algorithm computes the coefficients of  $P$  in the Bernstein bases of  $\ell, m$  and  $m, r$  from the coefficients of  $P$  in the Bernstein basis of  $\ell, r$ .

*Algorithm 10.2. [Bernstein Coefficients]*

- **Structure:** an archimedean real closed field  $\mathbb{R}$ .
- **Input:** a list  $b = b_0, \dots, b_p$  representing a polynomial  $P$  of degree  $\leq p$  in the Bernstein basis of  $\ell, r$ , and a number  $m \in \mathbb{R}$ .
- **Output:** the list  $b'$  representing  $P$  in the Bernstein basis of  $\ell, m$ , the list  $b''$  representing  $P$  in the Bernstein basis of  $m, r$ .
- **Complexity:**  $O(p^2)$ .
- **Procedure:**
  - Define  $\alpha = \frac{r-m}{r-\ell}, \beta = \frac{m-\ell}{r-\ell}$ .
  - Initialization:  $b_j^{(0)} := b_j, j = 0, \dots, p$ .
  - For  $i = 1, \dots, p$ ,
    - For  $j = 0, \dots, p-i$ , compute
 
$$b_j^{(i)} := \alpha b_j^{(i-1)} + \beta b_{j+1}^{(i-1)}$$
  - Output
 
$$b' = b_0^{(0)}, \dots, b_0^{(j)}, \dots, b_0^{(p)},$$

$$b'' = b_0^{(p)}, \dots, b_j^{(p-j)}, \dots, b_p^{(0)}.$$

Algorithm 10.2 (Bernstein Coefficients) can be visualized with the following triangle.



with  $b_j^{(i)} := \alpha b_j^{(i-1)} + \beta b_{j+1}^{(i-1)}, \alpha = (r-m)/(r-\ell), \beta = (m-\ell)/(r-\ell)$ .

The coefficients of  $P$  in the Bernstein basis of  $\ell, r$  appear in the top side of the triangle and the coefficients of  $P$  in the Bernstein basis of  $\ell, m$  and  $m, r$  appear in the two other sides of the triangle. Note that  $b_0^{(p)} = P(m)$ .

**Notation 10.36. [Reverted list]** We denote by  $\tilde{a}$  the list obtained by reversing the list  $a$ . □

**Proof of correctness:** It is enough to prove the part of the claim concerning  $\ell, m$ . Indeed, by Remark 10.31,  $\tilde{b}$  represents  $(-1)^p P$  in the Bernstein basis of  $r, \ell$ , and the claim is obtained by applying Algorithm 10.2 (Bernstein Coefficients) to  $\tilde{b}$  at  $m$ . The output is  $\tilde{b}''$  and  $\tilde{b}$  and the conclusion follows using again Remark 10.31.

Let  $\delta_{p,i}$  be the list of length  $p + 1$  consisting all zeroes except a 1 at the  $i + 1$ -th place. Note that  $\delta_{p,i}$  is the list of coefficients of  $\text{Bern}_{p,i}(\ell, t)$  in the Bernstein basis of  $\ell, r$ . We will prove that the coefficients of  $\text{Bern}_{p,i}(\ell, m)$  in the Bernstein basis of  $\ell, m$  coincide with the result of Algorithm 10.2 (Bernstein Coefficients) performed with input  $\delta_{p,i}$ . The correctness of Algorithm 10.2 (Bernstein Coefficients) for  $\ell, m$  then follows by linearity.

First notice that, since  $\alpha = (r - m)/(r - \ell), \beta = (m - \ell)/(r - \ell)$ ,

$$\begin{aligned} \frac{X - \ell}{r - \ell} &= \beta \frac{X - \ell}{m - \ell}, \\ \frac{r - X}{r - \ell} &= \alpha \frac{X - \ell}{m - \ell} + \frac{m - X}{m - \ell}. \end{aligned}$$

Thus

$$\begin{aligned} \left(\frac{X - \ell}{r - \ell}\right)^i &= \beta^i \left(\frac{X - \ell}{m - \ell}\right)^i \\ \left(\frac{r - X}{r - \ell}\right)^{p-i} &= \sum_{k=0}^{p-i} \binom{p-i}{k} \alpha^k \left(\frac{X - \ell}{m - \ell}\right)^k \left(\frac{m - X}{m - \ell}\right)^{p-i-k}. \end{aligned}$$

It follows that

$$\text{Bern}_{p,i}(\ell, r) = \binom{p}{i} \sum_{j=i}^p \binom{p-i}{j-i} \alpha^{j-i} \beta^i \left(\frac{X - \ell}{m - \ell}\right)^j \left(\frac{m - X}{m - \ell}\right)^{p-j}. \tag{10.2}$$

Since

$$\begin{aligned} \binom{p}{i} \binom{p-i}{j-i} &= \binom{j}{i} \binom{p}{j}, \\ \text{Bern}_{p,i}(\ell, r) &= \sum_{j=i}^p \binom{j}{i} \alpha^{j-i} \beta^i \binom{p}{j} \left(\frac{X - \ell}{m - \ell}\right)^j \left(\frac{m - X}{m - \ell}\right)^{p-j}. \end{aligned}$$

Finally,

$$\text{Bern}_{p,i}(\ell, r) = \sum_{j=i}^p \binom{j}{i} \alpha^{j-i} \beta^i \text{Bern}_{p,j}(\ell, m).$$

On the other hand, we prove by induction on  $p$  that Algorithm 10.2 (Bernstein Coefficients) with input  $\delta_{p,i}$  outputs the list  $\delta'_{p,i}$  starting with  $i$  zeroes and with  $(j + 1)$ -th element  $\binom{j}{i} \alpha^{j-i} \beta^i$  for  $j = i, \dots, p$ .

The result is clear for  $p = i = 0$ . If Algorithm 10.2 (Bernstein Coefficients) applied to  $\delta_{p-1,i-1}$  outputs  $\delta'_{p-1,i-1}$ , the equality

$$\binom{j}{i} \alpha^{j-i} \beta^i = \alpha \binom{j-1}{i} \alpha^{j-i-1} \beta^i + \beta \binom{j-1}{i-1} \alpha^{j-i} \beta^{i-1}$$

proves by induction on  $j$  that Algorithm 10.2 (Bernstein Coefficients) applied to  $\delta_{p,i}$  outputs  $\delta'_{p,i}$ . So the coefficients of  $\text{Bern}_{p,i}(\ell, r)$  in the Bernstein basis of  $\ell, m$  coincide with the output of Algorithm 10.2 (Bernstein Coefficients) with input  $\delta_{p,i}$ . □

**Corollary 10.37.** *Let  $b, b'$  and  $b''$  be the lists of coefficients of  $P$  in the Bernstein basis of  $\ell, r, \ell, m,$  and  $m, r$  respectively.*

- Algorithm 10.2 (Bernstein Coefficients) outputs  $b'$  and  $b''$  when applied to  $b$  with weights

$$\alpha = \frac{r - m}{r - \ell}, \beta = \frac{m - \ell}{r - \ell}.$$

- Algorithm 10.2 (Bernstein Coefficients) outputs  $b$  and  $\tilde{b}''$  when applied to  $b'$  with weights

$$\alpha' = \frac{m - r}{m - \ell}, \beta' = \frac{r - \ell}{m - \ell}.$$

- Algorithm 10.2 (Bernstein Coefficients) outputs  $\tilde{b}'$  and  $b$  when applied to  $b''$  with weights

$$\alpha'' = \frac{r - \ell}{r - m}, \beta'' = \frac{\ell - m}{r - m}.$$

**Complexity analysis of Algorithm 10.2:** The number of multiplications in the algorithm is  $2p(p + 1)/2$ , the number of additions is  $p(p + 1)/2$ . □

The following variant of Algorithm 10.2 (Bernstein Coefficients) is useful in the case of a polynomial with integer coefficients in the Bernstein basis of  $\ell, r$  since it avoids denominators.

*Algorithm 10.3. [Special Bernstein Coefficients]*

- **Structure:** an archimedean real closed field  $R$ .
- **Input:** a list  $b = b_0, \dots, b_p$  representing a polynomial  $P$  of degree  $\leq p$  in the Bernstein basis of  $\ell, r$ .
- **Output:** the list  $b'$  representing  $2^p P$  in the Bernstein basis of  $\ell, (\ell + r)/2$  and the list  $b''$  representing  $2^p P$  in the Bernstein basis of  $(\ell + r)/2, r$ .
- **Complexity:**  $O(p^2)$ .

• **Procedure:**

- Initialization:  $b_j^{(0)} := b_j$ , for  $j = 0, \dots, p$ .
- For  $i = 1, \dots, p$ ,
  - For  $j = 0, \dots, p - i$ , compute
 
$$b_j^{(i)} := b_j^{(i-1)} + b_{j+1}^{(i-1)}.$$
- Output
 
$$b' = 2^p b_0^{(0)}, \dots, 2^{p-j} b_0^{(j)}, \dots, b_0^{(p)},$$

$$b'' = b_0^{(p)}, \dots, 2^{p-j} b_j^{(p-j)}, \dots, 2^p b_p^{(0)}.$$

**Complexity analysis:** The number of additions in the algorithm is  $p(p+1)/2$ . The number of multiplications by 2 is  $p(p+1)$ . Note that if  $b \in \mathbb{Z}^{p+1}$ , then  $b' \in \mathbb{Z}^{p+1}$  and  $b'' \in \mathbb{Z}^{p+1}$ . If the bitsize of the  $b_i$  is bounded by  $\tau$ , the bitsizes of the  $b'_i$  and  $b''_i$  is bounded by  $p + \tau$ . □

**Proposition 10.38.** *Let  $P$  be a univariate polynomial of degree  $p$ . Let  $b$  be the Bernstein coefficients of  $P$  on  $(\ell, r)$  and  $b'$  the Bernstein coefficients of  $P$  on  $(\ell', r')$ . Denoting by  $c_i = \binom{p}{i} b_i$ ,  $Q = \sum_{i=0}^p c_i X^i$ ,  $c'_i = \binom{p}{i} b'_i$ , and  $Q' = \sum_{i=0}^p c'_i X^i$ , we have*

$$Q' = T_{-1}(\text{Rec}_p(\text{Co}_{r'-\ell'}(T_{\ell-\ell'}(\text{Co}_{1/(r-\ell)}(\text{Rec}_p(T_1(Q))))))).$$

**Proof:** Apply Proposition 10.27. □

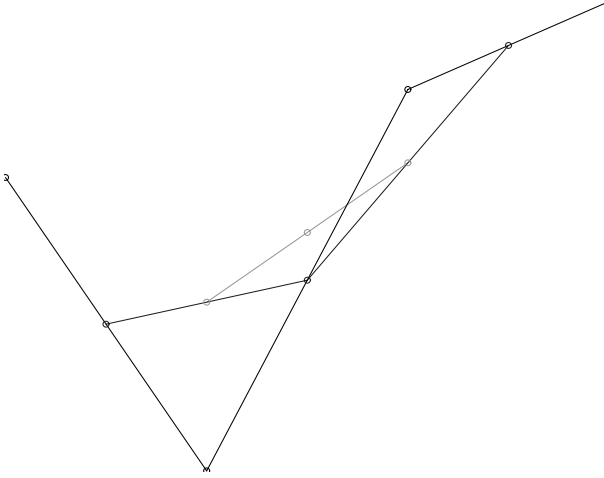
*Remark 10.39.* It is possible to output  $b'$  (and also  $b''$ ) with arithmetic complexity  $\tilde{O}(d)$  and binary complexity  $\tilde{O}(d(\tau + d))$  using Proposition 10.38 with  $\ell' = \ell$ ,  $r' = (\ell + r)/2$ , and Remark 8.7. □

Algorithm 10.2 (Bernstein Coefficients) gives a geometric construction of the control polygon of  $P$  on  $[\ell, m]$  and on  $[m, r]$  from the control polygon of  $P$  on  $[\ell, r]$ . The points of the new control polygons are constructed by taking iterated barycenters with weights  $\alpha$  and  $\beta$ . The construction is illustrated in Figure 10.3, where we show how the control line of  $P$  on  $[0, 1/2]$  is obtained from the control line of  $P$  on  $[0, 1]$ .

*Example 10.40.* Continuing Example 10.34, the Special Bernstein Coefficients Algorithm computes the following triangle.

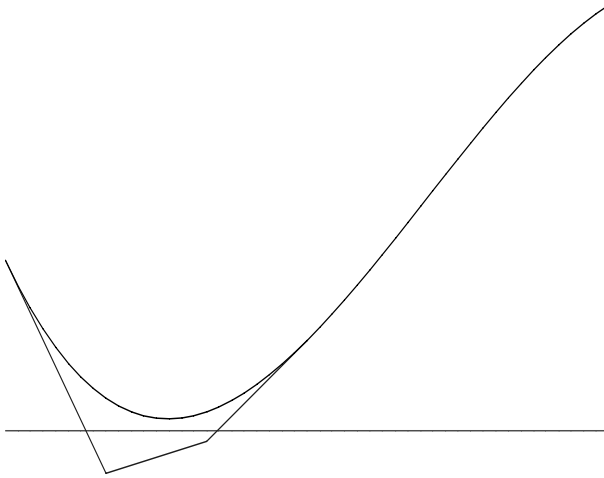
$$\begin{array}{cccc} 4 & -6 & 7 & 10 \\ & -2 & 1 & 17 \\ & & -1 & 18 \\ & & & 17 \end{array}$$

The Bernstein coefficients of  $8P$  on  $(0, 1/2)$  are 32,  $-8$ ,  $-2$ , 17, the Bernstein coefficients of  $8P$  on  $(1/2, 1)$  are 17, 36, 68, 80.



**Fig. 10.3.** Construction of the control line of  $P$  on  $[0, 1/2]$  by Bernstein Coefficients Algorithm

In Figure 10.4 we show the graph of  $P$  on  $[0, 1]$  and the control line on  $[0, 1/2]$ .



**Fig. 10.4.** Graph of  $P$  on  $[0, 1]$  and control line of  $P$  on  $[0, 1/2]$

□

We denote as usual by  $\text{Var}(b)$  the number of sign variations in a list  $b$ .

**Proposition 10.41.** *Let  $b, b'$  and  $b''$  be the lists of coefficients of  $P$  in the Bernstein basis of  $\ell, r, \ell, m$ , and  $m, r$ . If  $\ell < m < r$ , then*

$$\text{Var}(b') + \text{Var}(b'') \leq \text{Var}(b).$$

Moreover, if  $m$  is not a root of  $P$ ,  $\text{Var}(b) - \text{Var}(b') - \text{Var}(b'')$  is even.

**Proof:** The proof of the proposition is based on the following easy observations:

- Inserting in a list  $a = a_0, \dots, a_n$  a value  $x$  in  $[a_i, a_{i+1}]$  if  $a_{i+1} \geq a_i$  (resp. in  $[a_{i+1}, a_i]$  if  $a_{i+1} < a_i$ ) between  $a_i$  and  $a_{i+1}$  does not modify the number of sign variations.
- Removing from a list  $a = a_0, \dots, a_n$  with first non-zero  $a_k, k \geq 0$ , and last non-zero  $a_\ell, k \leq \ell \leq n$ , an element  $a_i, i \neq k, i \neq \ell$  decreases the number of sign variation by an even (possibly zero) natural number.

Indeed the lists

$$\begin{aligned} b &= b_0^{(0)}, \dots, \dots, \dots, \dots, b_p^{(0)} \\ b^{(1)} &= b_0^{(0)}, b_0^{(1)}, \dots, \dots, \dots, b_{p-1}^{(1)}, b_p^{(0)} \\ &\dots \\ b^{(i)} &= b_0^{(0)}, \dots, \dots, b_0^{(i)}, \dots, \dots, b_{p-i}^{(i)}, \dots, \dots, b_p^{(0)} \\ &\dots \\ b^{(p-1)} &= b_0^{(0)}, \dots, \dots, \dots, b_0^{(p-1)}, b_1^{(p-1)}, \dots, \dots, \dots, b_p^{(0)} \\ b^{(p)} &= b_0^{(0)}, \dots, \dots, \dots, \dots, b_0^{(p)}, \dots, \dots, \dots, b_p^{(0)} \end{aligned}$$

are successively obtained by inserting intermediate values and removing elements that are not end points, since when  $\ell < m < r$ ,  $b_j^{(i)}$  is between  $b_j^{(i-1)}$  and  $b_{j+1}^{(i-1)}$ , for  $i = 1, \dots, p, j = 0, \dots, p - i - 1$ . Thus  $\text{Var}(b^{(p)}) \leq \text{Var}(b)$  and the difference is even. Since

$$\begin{aligned} b' &= b_0^{(0)}, \dots, \dots, \dots, \dots, b_0^{(p)}, \\ b'' &= b_0^{(p)}, \dots, \dots, \dots, \dots, b_p^{(0)}, \end{aligned}$$

it is clear that

$$\text{Var}(b') + \text{Var}(b'') \leq \text{Var}(b^{(p)}) \leq \text{Var}(b).$$

If  $P(m) \neq 0$ , it is clear that

$$\text{Var}(b^{(p)}) = \text{Var}(b') + \text{Var}(b''), \text{ since } b_0^{(p)} = P(m) \neq 0. \quad \square$$

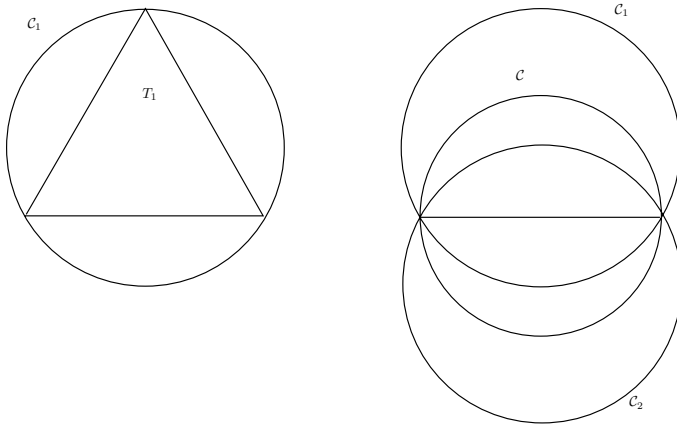
*Example 10.42.* Continuing Example 10.40, we observe, denoting by  $b, b'$  and  $b''$ , the lists of coefficients of  $P$  in the Bernstein basis of  $0, 1, 0, 1/2$ , and  $1/2, 1$ , that  $\text{Var}(b) = 2$ . This is visible in the figure: the control line for  $[0, 1]$  cuts twice the  $X$ -axis. Similarly,  $\text{Var}(b') = 2$ . This is visible in the figure: the control line for  $[0, 1/2]$  also cuts twice the  $X$ -axis. Similarly, it is easy to check that  $\text{Var}(b'') = 0$ .

We cannot decide from this information whether  $P$  has two roots in  $(0, 1/2)$  or no root in  $(0, 1/2)$ . □

Let  $b(P, (\ell, r))$  be the list of coefficients of  $P$  in the Bernstein basis of  $\ell, r$ ,  $\ell < r$ . The interval  $(\ell, r)$  is **active** if  $\text{Var}(b(P, (\ell, r))) > 0$ .

*Remark 10.43.* It is clear from Proposition 10.41 that if  $a_0 < \dots < a_N$ , the number of active intervals among  $(a_i, a_{i+1})$  is at most  $p$ . □

Let  $P \in \mathbb{R}[X]$  and let  $b(P, (\ell, r))$  be the list of coefficients of  $P$  in the Bernstein basis of  $\ell, r$ . We describe cases where the number  $\text{Var}(b(P, (\ell, r)))$  coincides with the number of roots of  $P$  on  $(\ell, r)$ . Denote by  $\mathcal{C}(\ell, r)$  the closed disk with  $[\ell, r]$  as a diameter, by  $\mathcal{C}_1(\ell, r)$  the closed disk whose boundary circumscribes the equilateral triangle  $T_1$  based on  $[\ell, r]$  (see Figure 10.5), and by  $\mathcal{C}_2(\ell, r)$  the closed disk symmetric to  $\mathcal{C}_1(\ell, r)$  with respect to the  $X$ -axis (see Figure 10.5)



**Fig. 10.5.**  $\mathcal{C}(\ell, r)$ ,  $\mathcal{C}_1(\ell, r)$  and  $\mathcal{C}_2(\ell, r)$

**Theorem 10.44. [Theorem of three circles]**

- Let  $P$  be a separable polynomial of  $\mathbb{R}[X]$ .
- If  $P$  has no root in  $\mathcal{C}(\ell, r)$ , then  $\text{Var}(b(P, (\ell, r))) = 0$ .
- If  $P$  has exactly one root in  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$ , then  $\text{Var}(b(P, (\ell, r))) = 1$ .

**Proof:** We identify  $\mathbb{R}^2$  with  $\mathbb{C} = \mathbb{R}[i]$ . The image of the complement of  $\mathcal{C}(\ell, r)$  (resp.  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$ ) under translation by  $-\ell$  followed by contraction by ratio  $r - \ell$  is the complement of  $\mathcal{C}(0, 1)$  (resp.  $\mathcal{C}_1(0, 1) \cup \mathcal{C}_2(0, 1)$ ). The image of the complement of  $\mathcal{C}(0, 1)$  under the inversion  $z \mapsto 1/z$  is the half plane of complex numbers with real part less than 1. Translating by  $-1$ , we get the half plane of complex numbers with non-positive real part.



The image of the complement of  $\mathcal{C}_1(0, 1) \cup \mathcal{C}_2(0, 1)$ , under the inversion  $z \mapsto 1/z$  is the sector

$$\{(x + iy) \in \mathbb{R}[i] \mid |y| \leq \sqrt{3}(1 - x)\}.$$

Translating this region by  $-1$ , we get the cone  $\mathcal{B}$  defined in Proposition 2.40.

The statement follows from Proposition 2.39, Proposition 2.44 and Proposition 10.27.  $\square$

**Corollary 10.45.** *If  $P$  is separable,  $\text{Var}(b(P, (\ell, r))) \geq 2$  implies that  $P$  has at least two roots in  $\mathcal{C}(\ell, r)$  or the interval  $(\ell, r)$  contains exactly one real root and  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$  contains a pair of conjugate roots.*

**Proof:** If  $P$  has no root in  $\mathcal{C}(\ell, r)$ , then  $\text{Var}(b(P, (\ell, r))) = 0$ , by Theorem 10.44. Thus,  $P$  has at least one root in  $\mathcal{C}(\ell, r)$ . If this is the only root in  $\mathcal{C}(\ell, r)$ , the root is in  $(\ell, r)$  and  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$  must contain a pair of conjugate roots because otherwise  $\text{Var}(b(P, (\ell, r))) = 1$ , by Theorem 10.44.  $\square$

Suppose that  $P \in \mathbb{R}[X]$  is a polynomial of degree  $p$  with all its real zeroes in  $(-2^N, 2^N)$  (where  $N$  is a natural number) and let  $\bar{P}$  be the separable part of  $P$ . Consider natural numbers  $k$  and  $c$  such that  $0 \leq c \leq 2^k$  and define

$$\begin{aligned} \ell &= \frac{-2^{N+k} + c2^{N+1}}{2^k} \\ r &= \frac{-2^{N+k} + (c+1)2^{N+1}}{2^k} \end{aligned}$$

Let  $b(\bar{P}, \ell, r)$  denote the list of coefficients of  $2^{kp}\bar{P}$  in the Bernstein basis of  $(\ell, r)$ . Note that if  $\bar{P}$  is such that its list of coefficients in the Bernstein basis of  $(-2^N, 2^N)$  belong to  $\mathbb{Z}$ , the coefficients of  $2^{kp}\bar{P}$  in the Bernstein basis of  $(\ell, r)$  belong to  $\mathbb{Z}$ . This follows clearly from Algorithm 10.3 (Special Bernstein Coefficients).

*Remark 10.46.* Let  $\text{sep}$  be the minimal distance between two roots of  $P$  in  $\mathbb{C}$ , and let  $N$  be such that the real roots of  $P$  belong to  $(-2^N, 2^N)$ , and  $k \geq -\log_2(\text{sep}) + N + 1$ . Since the circle of center  $(\ell+r)/2, 0$  and radius  $r - \ell$  contains  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$ , and two points inside this circle have distance at most  $2(r - \ell)$ , it is clear that the polynomial  $\bar{P}$  has at most one root in  $(\ell, r)$  and has no other complex root in  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$ . So,  $\text{Var}(b(\bar{P}, \ell, r))$  is zero or one, using Theorem 10.44.

Thus, it is possible to decide, whether  $\bar{P}$  has exactly one root in  $(\ell, r)$  or has no root on  $(\ell, r)$ , by testing whether  $\text{Var}(b(\bar{P}, \ell, r))$  is zero or one.  $\square$

*Example 10.47.* Continuing Example 10.42, let us study the roots of  $P$  on  $[0, 1]$ , as a preparation to a more formal description of Algorithm 10.4 (Real Root Isolation).

The Bernstein coefficients of  $P$  for  $0, 1$  are  $4, -6, 7, 10$ . There maybe roots of  $P$  on  $(0, 1)$  as there are sign variations in these Bernstein coefficients.

As already seen in Example 10.42, a first application of Algorithm 10.3 (Special Bernstein Coefficients) gives

$$\begin{array}{cccc} 4 & -6 & 7 & 10 \\ & -2 & 1 & 17 \\ & & -1 & 18 \\ & & & 17 \end{array}$$

There maybe roots of  $P$  on  $(0, 1/2)$  as there are sign variations in the Bernstein coefficients of  $8P$  on  $(0, 1/2)$  which are  $32, -8, -2, 17$ . There are no roots of  $P$  on  $(1/2, 1)$  as there are no sign variations in the Bernstein coefficients of  $8P$  on  $(1/2, 1)$  which are  $17, 36, 68, 80$ .

Let us apply once more Algorithm 10.3 (Special Bernstein Coefficients):

$$\begin{array}{cccc} 32 & -8 & -2 & 17 \\ & 24 & -10 & 15 \\ & & 14 & 5 \\ & & & 19 \end{array}$$

The Bernstein coefficients of  $64P$  on  $(0, 1/4)$  are  $256, 96, 28, 19$ , and the Bernstein coefficients of  $64P$  on  $(1/4, 1/2)$  are  $19, 10, 60, 136$ . There are no sign variations on the sides of the triangle so there are no roots of  $P$  on  $(0, 1/4)$  and on  $(1/4, 1/2)$ .

Finally there are no roots of  $P$  on  $[0, 1]$ . □

**Definition 10.48. [Isolating list]** Let  $Z$  be a finite subset of  $\mathbb{R}$ . An **isolating list for  $Z$**  is a finite list  $L$  of rational points and open intervals with rational end points of  $\mathbb{R}$ , such that each element of  $L$  contains exactly one element of  $Z$ , every element of  $Z$  belongs to an element of  $L$  and two elements of  $L$  have an empty intersection. □

*Algorithm 10.4. [Real Root Isolation]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:** a non-zero polynomial  $P \in \mathbb{Z}[X]$ .
- **Output:** a list isolating for the zeroes of  $P$  in  $\mathbb{R}$ .
- **Binary complexity:**  $O(p^5(\tau + \log_2(p))^2)$ , where  $p$  is a bound on the degree of  $P$ , and  $\tau$  is a bound on the bitsize of the coefficients of  $P$ .
- **Procedure:**
  - Compute  $N \in \mathbb{N}$ ,  $N \geq \log_2(C(P))$  (Notation 10.1) such that  $(-2^N, 2^N)$  contains the roots of  $P$  in  $\mathbb{R}$ .
  - Compute  $\bar{P}$ , the separable part of  $P$  using Algorithm 10.1 (Gcd and Gcd-Free part). Replace  $\bar{P}$  by  $\lambda(\bar{P}, -2^N, 2^N) \bar{P}$ , using Corollary 10.30 and its notation. Compute  $b(\bar{P}, -2^N, 2^N)$ , the Bernstein coefficients of  $\bar{P}$  on  $(-2^N, 2^N)$ , using Remark 10.28.

- Initialization:  $\text{Pos} := \{b(\bar{P}, -2^N, 2^N)\}$  and  $L(P)$  is the empty list.
- While  $\text{Pos}$  is non-empty,
  - Remove an element  $b(\bar{P}, \ell, r)$  from  $\text{Pos}$ .
  - If  $\text{Var}(b(\bar{P}, \ell, r)) = 1$ , add  $(\ell, r)$  to  $L(P)$ .
  - If  $\text{Var}(b(\bar{P}, \ell, r)) > 1$ ,
    - Compute  $b(\bar{P}, \ell, m)$  and  $b(\bar{P}, m, r)$ , with  $m = (\ell + r)/2$ , using Algorithm 10.3 (Special Bernstein Coefficients) and add them to  $\text{Pos}$ .
    - If the sign of  $\bar{P}(m)$  is 0, see Remark 10.31 add  $\{m\}$  to  $L(P)$ .

**Proof of correctness:** The algorithm terminates since  $\mathbb{R}$  is archimedean, using Remark 10.46. Its correctness follows from Theorem 10.44. Note that, since there is only one root of  $\bar{P}$  on each interval  $[a, b]$  of  $L(P)$ , we have  $\bar{P}(a)\bar{P}(b) < 0$ .  $\square$

The complexity analysis requires some preliminary work.

*Remark 10.49.* Note that by Corollary 10.45, the binary tree  $T$  produced by Algorithm 10.4 enjoys the following properties:

- the interval labeling the root of the tree  $T$  contains all roots of  $P$  in  $\mathbb{R}$ ,
- at every leaf node labelled by  $(\ell, r)$  of  $T$ , the interval  $(\ell, r)$  contains either no root or one single root of  $P$ ,
- at every node labelled by  $(\ell, r)$  of  $T$  which is not a leaf, either  $P$  has at least two roots in  $\mathcal{C}(\ell, r)$ , or the interval  $(\ell, r)$  contains exactly one real root and the union of the two circles  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$  contains two conjugate roots.  $\square$

So, we consider binary trees labeled by open intervals with rational endpoints, such that if a node of the tree is labeled by  $(\ell, r)$ , its children are labeled either by  $(\ell, m)$  or by  $(m, r)$ , with  $m = (\ell + r)/2$ . Such a tree  $T$  is an **isolating tree** for  $P$  if the following properties holds:

- the interval labeling the root of the tree  $T$  contains all the roots of  $P$  in  $\mathbb{R}$ ,
- at every leaf node labelled by  $(\ell, r)$  of  $T$ , the interval  $(\ell, r)$  contains either no root or one single root of  $P$ ,
- at every node labelled by  $(\ell, r)$  of  $T$  which is not a leaf, either  $P$  has at least two roots in  $\mathcal{C}(\ell, r)$ , or the interval  $(\ell, r)$  contains exactly one root of  $P$  and the union of the two circles  $\mathcal{C}_1(\ell, r) \cup \mathcal{C}_2(\ell, r)$  contains two conjugate roots.

As noted above, the binary tree produced by Algorithm 10.4 (Real Root Isolation) is an isolating tree for  $P$ .

Let  $P \in \mathbb{Z}[X]$  be of degree  $p$ , and let  $\tau$  be a bound on the bitsize of the coefficients of  $P$  and  $\nu$  a bound on the bitsize of  $p$ . By Corollary 10.4 all the roots of  $P$  belong to the interval

$$u_0 = (-2^{\tau+\nu}, 2^{\tau+\nu}). \quad (10.3)$$

**Proposition 10.50.** *Let  $T$  be an isolating tree for  $P$  with root  $u_0$  and  $L$  its set of leaves. Given a leaf  $u \in L$ , denote by  $h_u$  its depth. Then,*

$$\sum_{u \in L} h_u \leq 2(2\tau + 3\nu + 3)p.$$

Before proving Proposition 10.50 we need to study in some properties of  $T$  in more detail. Note that a node of  $T$  is labeled by an interval  $(\ell, r)$ . Note also that a leaf of  $T$  is

- either a leaf of type 1, when  $P$  has a root on  $(\ell, r)$ ,
- or a leaf of type 0, when  $P$  has no root on  $(\ell, r)$ .

In order to bound the number of nodes of  $T$  we introduce a subtree  $T'$  of  $T$  defined by pruning certain leaves from  $T$ :

- If a leaf  $u$  has a sibling that is not a leaf, we prune  $u$ .
- If  $u$  and  $v$  are both leaves and siblings of each other, then we prune exactly one of them; the only constraint is that a leaf of type 0 is pruned preferably to a leaf of type 1.

We denote by  $L'$  the set of leaves of  $T'$ .

Clearly,

$$\sum_{u \in L} h_u \leq 2 \sum_{u \in L'} h_u \tag{10.4}$$

So in order to bound  $\sum_{u \in L} h_u$  it suffices to bound  $\sum_{u \in L'} h_u$ .

If  $u = (\ell, r)$  is an interval, we denote by  $w(u) = r - \ell$  the **width** of the interval  $u$ . We define  $w_0 = w(u_0)$  where  $u_0$  is the interval labeling the root of the tree  $T$ . The number of nodes along the path from any  $u \in L'$  to the root of  $T'$  is exactly  $\log_2(w_0/w(u))$ . Thus

$$\sum_{u \in L'} h_u \leq \sum_{u \in U} \log_2\left(\frac{w_0}{w(u)}\right). \tag{10.5}$$

Let  $u \in L'$  be a root of  $T'$ . We are going to define two roots of  $P$ ,  $\alpha_u$  and  $\beta_u$  of  $P$  such that

$$|\alpha_u - \beta_u| < 4w(u).$$

Furthermore we will show that if  $u, u'$  have the same type (both type 0, or both type 1) then  $\{\alpha_u, \beta_u\}$  and  $\{\alpha_{u'}, \beta_{u'}\}$  are disjoint.

Let  $v$  be the parent of the leaf  $u$ .

- a) If  $u$  is of type 1, then  $u$  contains a root  $\alpha_u$ , and the union of the two circles  $\mathcal{C}_1(v) \cup \mathcal{C}_2(v)$  contains two conjugate roots, and we denote by  $\beta_v$  one of these. Then

$$|\alpha_u - \beta_u| < (2/\sqrt{3})w(v) = (4/\sqrt{3})w(u) < 4w(u). \tag{10.6}$$

Let  $u'$  be another leaf of type 1 and  $v'$  its parent. Clearly,  $\alpha_u \neq \alpha_{u'}$ . We claim that it is possible to choose  $\beta_u$  and  $\beta_{u'}$  such that  $\beta_u \neq \beta_{u'}$ . Consider the case when  $v$  and  $v'$  are siblings. Moreover, assume that  $\beta_u$  and  $\overline{\beta_u}$  are the only non-real roots in  $\mathcal{C}_1(v) \cup \mathcal{C}_2(v)$  and  $\mathcal{C}_1(v') \cup \mathcal{C}_2(v')$ . Then it must be the case that either  $\beta_u \in \mathcal{C}_1(v) \cap \mathcal{C}_1(v')$ , or  $\beta_u \in \mathcal{C}_2(v) \cap \mathcal{C}_2(v')$ . In either case, we can choose  $\beta_{u'} = \overline{\beta_u}$ , distinct from  $\beta_u$ . Thus  $\{\alpha_u, \beta_u\}$  and  $\{\alpha_{u'}, \beta_{u'}\}$  are disjoint.

- b) If  $u$  is of type 0,  $P$  has one root  $\alpha_u$  in  $\mathcal{C}(v)$ . Clearly,  $\alpha_u$  is non-real, otherwise  $u$  would either have a non-leaf or a leaf of type 1 as a sibling in  $T$ , and would have been pruned from  $T$ . Thus  $\mathcal{C}(v)$  contains  $\overline{\alpha_u} \neq \alpha_u$  and we define  $\beta_u = \overline{\alpha_u}$ . Then,

$$|\alpha_u - \beta_u| < 2w(u) < 4w(u). \tag{10.7}$$

If  $u'$  is another leaf of type 1, then  $\{\alpha_u, \overline{\alpha_u}\}$  and  $\{\alpha_{u'}, \overline{\alpha_{u'}}\}$  are disjoint, since  $\mathcal{C}(v)$  and  $\mathcal{C}(v')$  are disjoint.

Taking logarithms and substituting (10.6) and (10.7) in (10.5) we get

$$\sum_{u \in L'} h_u \leq \sum_{u \in U} \log_2 \left( \frac{4w_0}{|\alpha_u - \beta_u|} \right) \tag{10.8}$$

**Lemma 10.51.** *We have  $\#(L') \leq p$ . More precisely denoting by  $L_0$  the leaves of type 0 of  $T'$  and by  $L_1$  the leaves of type 1 of  $T'$ ,*

- a)  $\#(L_0)$  is at most half the number of non-real roots of  $P$ .
- b)  $\#(L_1)$  is at most the number of real roots of  $P$ .

**Proof:** As shown above, to every  $u \in L_0$  we can associate a unique pair of non-real roots  $(\alpha_u, \beta_u)$ ,  $\beta_u = \overline{\alpha_u}$ . Since the non-real roots come in pair, the upper bound of  $U_0$  follows.

Again by the arguments above, to each  $u \in L_1$  we can associate a unique real root  $\alpha_u$  and the claim on  $L_1$  follows.

Finally  $\#(L') \leq p$ . □

**Proof of Proposition 10.50:**

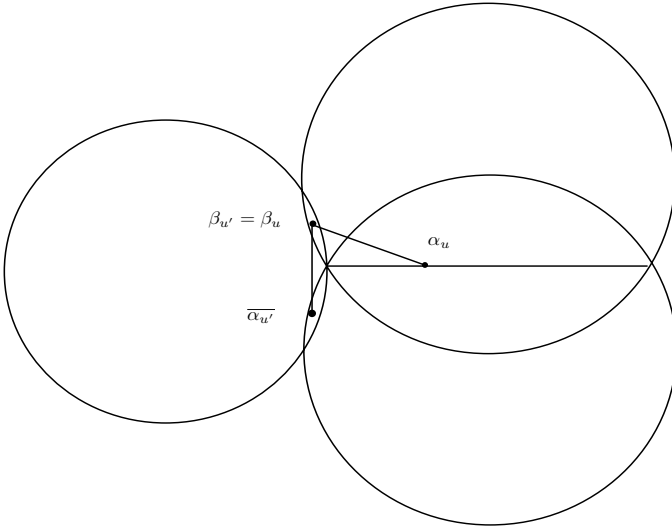
From Corollary 10.4, we know that  $\log_2(w_0) \leq \tau + \nu + 1$ , where  $\nu$  is a bound on the bitsize of  $p$ , and from Lemma 10.51, we have  $\#(L') \leq p$ . So, from (10.8), we obtain

$$\sum_{u \in L'} h_u \leq (\tau + \nu + 3)p - \sum_{u \in L'} \log_2(|\alpha_u - \beta_u|). \tag{10.9}$$

It remains to lower bound  $\sum_{u \in L'} \log_2(|\alpha_u - \beta_u|)$ . This will be done using Corollary 10.24. Consider the graph  $G$  whose edge set is  $E_0 \cup E_1$  where

$$\begin{aligned} E_0 &= \{(\alpha_u, \beta_u) \mid u \in L_0\}, \\ E_1 &= \{(\alpha_u, \beta_u) \mid u \in L_1\}. \end{aligned}$$

We want to show that  $G$  satisfies the hypotheses of Corollary 10.24. First of all, for any  $u \in L'$ , we can reorder the pair  $(\alpha_u, \beta_u)$  such that  $|\alpha_u| \leq |\beta_u|$ , without affecting (10.9).



**Fig. 10.6.** A type 0 leaf and a type 1 leaf sharing the same root

Now we show that the in-degree of  $G$  is at most 1. Clearly the edge sets  $E_0$  and  $E_1$  have in-degree at most 1. However in  $E_0 \cup E_1$ , a case like that illustrated in Figure 10.6 can happen. That is, a  $u \in L_1$  and a  $u' \in L_0$  such that  $\beta_u = \beta_{u'}$ . But in such a case we can always reorder the edge  $(\alpha_{u'}, \beta_{u'})$  to  $(\beta_{u'}, \alpha_{u'})$  since  $\beta_{u'} = \overline{\alpha_{u'}}$ , and thus reduce the in-degree to 1.

Now we may apply Corollary 10.24 to  $G$ , and get

$$\prod_{u \in L'} |\alpha_u - \beta_u| \geq (p/\sqrt{3})^{-\#(U)} p^{-p/2} (p+1)^{(1-p)/2} 2^{\tau(1-p)}.$$

Taking logarithms of both sides gives

$$-\sum_{u \in L'} \log_2(|\alpha_u - \beta_u|) \leq (\tau + 2\nu) p, \tag{10.10}$$

since  $\#(L') \leq p$  by Lemma 10.51, and  $\log_2(p) < \log_2(p+1) \leq \nu$ .

Substituting in (10.9), we obtain

$$\sum_{u \in L'} h_u \leq (\tau + \nu + 3) p + (\tau + 2\nu) p = (2\tau + 3\nu + 3) p.$$

The proposition is proved, since  $\sum_{u \in L} h_u \leq 2 \sum_{u \in L'} h_u$ . □

**Complexity analysis of Algorithm 10.4:** The computation of  $\bar{P}$  takes  $O(p^2)$  arithmetic operations using Algorithm 10.1.

Using Proposition 10.50 and Remark 10.49, the number of nodes produced by the algorithm is at most  $O(p(\tau + \log_2(d)))$ .

At each node the computation takes  $O(p^2)$  arithmetic operations using Algorithm 10.3.

It follows from Corollary 10.12 and Lemma 10.2 that the coefficients of  $\bar{P}$  are of bitsize  $O(\tau + p)$ .

The coefficients of  $b(\bar{P}, \ell, m)$  and  $b(\bar{P}, m, r)$  are integer numbers of bitsizes  $O(p^2(\tau + \log_2(p)))$  according to Corollary 10.30.

Since there are only additions and multiplications by 2 to perform, the estimate for the binary complexity of the algorithm is  $O(p^5(\tau + \log_2(p))^2)$ .  $\square$

*Remark 10.52.* Using Remark 10.39, and the preceding complexity analysis, it is possible to compute the output of Algorithm 10.4 (Real Root Isolation) with complexity  $\tilde{O}(p^2 \tau)$  and binary complexity  $\tilde{O}(p^4 \tau^2)$ . The same remark applies for other algorithms in this section: it is possible to compute the output of Algorithm 10.6 and Algorithm 10.7 with complexity  $\tilde{O}(p^2 \tau)$  and binary complexity  $\tilde{O}(p^4 \tau^2)$  and the output of Algorithm 10.8 with complexity  $\tilde{O}(s^2 p^2 \tau)$  and binary complexity  $\tilde{O}(s^2 p^4 \tau^2)$ .  $\square$

*Remark 10.53.* It is clear that Algorithm 10.4 (Real Root Isolation) also provides a method for counting real roots.

Note that if  $(\ell, r)$  is an open interval isolating a root  $x$  of  $P$ , and  $m = (\ell + r)/2$  it is easy to decide whether  $x = m$ ,  $x \in (\ell, m)$ , or  $x \in (m, r)$  by applying once more Algorithm 10.3 (Special Bernstein Coefficients), since the sign of  $P(m)$  is part of the output.  $\square$

We now give a variant of Algorithm 10.4 (Real Root Isolation) where the computations take place in the basis of monomials.

**Notation 10.54.** Let  $P \in \mathbb{Z}[X]$  be a polynomial of degree  $p$  having all its real roots between  $(-2^N, 2^N)$ . We define for  $r - \ell = a \cdot 2^k$ ,  $a \in \mathbb{Z}$ ,  $k \in \mathbb{Z}$ ,  $a$  odd,

$$\begin{aligned} P[\ell, r] &= P((r - \ell)X + \ell) && \text{if } k \geq 0 \\ &= 2^{kp} P((r - \ell)X + \ell) && \text{if } k < 0 \end{aligned}$$

In other words, the roots of  $P[\ell, r]$  in  $(0, 1)$  are in 1-1 correspondence with the roots of  $P$  on  $(\ell, r)$ , and  $P[\ell, r] \in \mathbb{Z}[X]$ . Note that

$$\begin{aligned} P[-2^N, 2^N] &= \text{Co}_{2^{N+1}}(\text{T}_{2^N}(P)) \\ &= \text{Co}_2(\text{T}_{-1}(\text{Co}_{2^N}(P))) \end{aligned} \tag{10.11}$$

and, if  $m = (\ell + r)/2$ .

$$\begin{aligned} P[\ell, m] &= 2^q \text{Co}_{1/2}(P[\ell, r]) && \text{if } m \notin \mathbb{Z}, \\ &= \text{Co}_{1/2}(P[\ell, r]) && \text{if } m \in \mathbb{Z}, \\ P[m, r] &= \text{T}_{-1}(P[\ell, m]). \end{aligned} \tag{10.12}$$

$\square$

The following proposition explains how to recover the sign variations in the Bernstein coefficients of  $P$  for  $(\ell, r)$ , from  $P[\ell, r]$ .

**Proposition 10.55.**

$$\text{Var}(b(P, (\ell, r))) = \text{Var}(\text{T}_{-1}(\text{Rec}_p(P[\ell, r])))$$

**Proof:** Follows immediately from Proposition 10.27. □

*Algorithm 10.5.* **[Descartes' Real Root Isolation]**

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:** a non-zero polynomial  $P \in \mathbb{Z}[X]$ .
- **Output:** a list isolating for the zeroes of  $P$  in  $\mathbb{R}$ .
- **Binary complexity:**  $O(p^5(\tau + \log_2(p))^2)$ , where  $p$  is a bound on the degree of  $P$ , and  $\tau$  is a bound on the bitsize of the coefficients of  $P$ .
- **Procedure:**
  - Compute  $N \in \mathbb{N}$ ,  $N \geq \log_2(C(P))$  (Notation 10.1) such that  $(-2^N, 2^N)$  contains the roots of  $P$  in  $\mathbb{R}$ .
  - Compute  $\bar{P}$ , the separable part of  $P$  using Algorithm 10.1 (Gcd and Gcd-Free part) and denote by  $q$  its degree.
  - Compute  $\bar{P}[-2^N, 2^N] = \text{Co}_2(\text{T}_{-1}(\text{Co}_{2^N}(\bar{P})))$ , using Algorithm 8.10.
  - Initialization:  $\text{Pos} := \{\bar{P}[-2^N, 2^N]\}$  and  $L(P)$  is the empty list.
  - While  $\text{Pos}$  is non-empty,
    - Remove an element  $\bar{P}[\ell, r]$  from  $\text{Pos}$ .
    - If  $\text{Var}(\text{T}_{-1}(\text{Rec}_q(\bar{P}[\ell, r]))) = 1$ , add  $(\ell, r)$  to  $L(P)$ .
    - If  $\text{Var}(\text{T}_{-1}(\text{Rec}_q(\bar{P}[\ell, r]))) > 1$ ,
      - Let  $m = (\ell + r)/2$ . Compute,

$$\begin{aligned} \bar{P}[\ell, m] &= 2^q \text{Co}_{1/2}(\bar{P}[\ell, r]) \text{ if } m \notin \mathbb{Z}, \\ &= \text{Co}_{1/2}(\bar{P}[\ell, r]) \text{ if } m \in \mathbb{Z}, \\ \bar{P}[m, r] &= \text{T}_{-1}(\bar{P}[\ell, m]), \end{aligned}$$

using Algorithm 8.10, add

$$\bar{P}[\ell, m] \text{ and } \bar{P}[m, r]$$

to  $\text{Pos}$ .

- If the sign of  $\bar{P}(m)$  is 0, add  $\{m\}$  to  $L(P)$ .

**Proof of correctness:** The algorithm terminates since  $\mathbb{R}$  is archimedean. The correctness of the algorithm follows from the correctness of Algorithm 10.4 (Real Root Isolation), using Proposition 10.55. □



**Complexity analysis:** The computation of  $\bar{P}$  takes  $O(p^2)$  arithmetic operations using Algorithm 10.1. Using Proposition 10.50 and Remark 10.49, the number of nodes produced by the algorithm is at most  $O(p(\tau + \log_2(d)))$ . At each node the computation takes  $O(p^2)$  arithmetic operations. It follows from Corollary 10.12 and Lemma 10.2 that the coefficients of  $\bar{P}$  are of bitsize  $O(\tau + p)$ . Finally the coefficients of  $\bar{P}[\ell, r]$  and  $T_{-1}(\text{Rec}(\bar{P}[\ell, r]))$  are bounded by  $O(p^2(\tau + \log_2(p)))$  and only multiplication by 2 and additions are performed.  $\square$

To evaluate the sign of another polynomial  $Q$  at the root of a polynomial characterized by an isolating interval, it may be necessary to refine the isolating intervals further. We need the following definition.

**Definition 10.56. [Isolating list with signs]** Let  $Z$  be a finite subset of  $\mathbb{R}$  and a finite list  $\mathcal{Q}$  of polynomial of  $\mathbb{R}[X]$ . An **isolating list with signs for  $Z$  and  $\mathcal{Q}$**  is a finite list  $L$  of couples  $(I, \sigma)$  such that  $I$  is a rational point or an open interval with rational end points, and  $\sigma$  is an element of  $\{-1, 0, 1\}^{\mathcal{Q}}$ . Every element of  $Z$  belongs to some  $I$  with  $(I, \sigma)$  in  $L$  and for every  $(I, \sigma)$  in  $L$ , there exists one and only one element  $x$  in  $I$  and  $\sigma$  is the sign condition realized by the family  $\mathcal{Q}$  at  $x$ .  $\square$

*Algorithm 10.6. [Sign at a Real Root]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:** a polynomial  $P \in \mathbb{Z}[X]$ , a list  $L(P)$  isolating for the zeroes of  $P$  in  $\mathbb{R}$  and a polynomial  $Q \in \mathbb{Z}[X]$ .
- **Output:** an isolating list with signs for the zeroes of  $P$  in  $\mathbb{R}$  and  $\{Q\}$ .
- **Binary complexity:**  $O(p^5(\tau + \log_2(p))^2)$ , where  $p$  is a bound on the degree of  $P$ , and  $\tau$  is a bound on the bitsize of the coefficients of  $P$ .
- **Procedure:**
  - First step: Identify the common roots of  $P$  and  $Q$  as follows. This is done as follows:
  - Compute the greatest common divisor  $G$  of  $\bar{P}$  and  $Q$ . Note that  $G$  is separable. If the structure is  $\mathbb{Z}$ , replace  $G$  by  $\lambda(G, -2^N, 2^N)G$  using Corollary 10.30 and its notation.
  - Initialization:
    - Set  $N(P) := \emptyset$ ,  $\text{NCom}(P, Q) := \emptyset$  ( $\text{NCom}(P, Q)$  will contain points or intervals corresponding to roots of  $P$  which are not roots of  $Q$ ). For every  $\{a\} \in L(P)$ , add  $(\{a\}, \text{sign}(Q(a)))$  to  $N(P)$ .
    - Compute  $b(G, \ell, r)$ , the Bernstein coefficients of  $G$ , for the intervals  $(\ell, r) \in L(P)$  using Proposition 10.27. Set
$$\text{Pos} := \{(b(\bar{P}, \ell, r), b(G, \ell, r))\}$$
 for the intervals  $(\ell, r) \in L(P)$ .
  - While  $\text{Pos}$  is non-empty,
    - Remove an element  $b(\bar{P}, \ell, r), b(G, \ell, r)$  from  $\text{Pos}$ .
    - If  $\text{Var}(b(G, \ell, r)) = 1$ , add  $((\ell, r), 0)$  to  $N(P)$ .

- If  $\text{Var}(b(G, \ell, r)) = 0$ , add  $(b(\bar{P}, (\ell, r)))$  to  $\text{NCom}(P, Q)$ .
- If  $\text{Var}(b(G, \ell, r)) > 1$ , compute  $(b(\bar{P}, \ell, m), b(G, \ell, m))$  and  $(b(\bar{P}, m, r), b(G, m, r))$  with  $m = (\ell + r)/2$  using Algorithm 10.3 (Special Bernstein Coefficients).
  - If  $\bar{P}(m) = 0$ , add  $(\{m\}, \text{sign}(Q(m)))$  to  $N(P)$ .
  - If the signs of  $\bar{P}$  at the right of  $\ell$  and at  $m$  coincide, add  $(b(\bar{P}, m, r), b(G, m, r))$  to  $\text{Pos}$ .
  - If the signs of  $\bar{P}$  at the right of  $\ell$  and at  $m$  differ, add  $(b(\bar{P}, \ell, m), b(G, \ell, m))$  to  $\text{Pos}$ .
- Second step: Find the sign of  $Q$  at the roots of  $P$  where  $Q$  is non-zero. This is done as follows:
  - Initialization:  $\text{Pos} := \text{NCom}(P, Q)$ . If the structure is  $\mathbb{Z}$ , replace  $Q$  by  $\lambda(Q, -2^N, 2^N)Q$  using Corollary 10.30 and its notation.
  - While  $\text{Pos}$  is non-empty,
    - Remove an element  $b(\bar{P}, \ell, r)$  from  $\text{Pos}$ . Compute  $b(Q, \ell, r)$  the Bernstein coefficients of  $Q$  on  $(\ell, r)$  using Proposition 10.27.
    - If  $\text{Var}(b(Q, \ell, r)) = 0$ , add  $((\ell, r), \tau)$  to  $N(P)$ , where  $\tau$  is the sign of any element of  $b(Q, \ell, r)$ .
    - If  $\text{Var}(b(Q, \ell, r)) \neq 0$ , compute  $b(\bar{P}, \ell, m)$  and  $b(\bar{P}, m, r)$  using Algorithm 10.3 (Special Bernstein Coefficients).
      - If  $\bar{P}(m) = 0$ , add  $(\{m\}, \text{sign}(Q(m)))$  to  $N(P)$ .
      - If the signs of  $\bar{P}$  at the right of  $\ell$  and at  $m$  coincide, add  $(b(\bar{P}, m, r))$  to  $\text{Pos}$ .
      - If the signs of  $\bar{P}$  at the right of  $\ell$  and at  $m$  differ, add  $(b(\bar{P}, \ell, m))$  to  $\text{Pos}$ .

**Proof of correctness:** The algorithm terminates since  $\mathbb{R}$  is archimedean. Its correctness follows from Theorem 10.37. Note that on any interval output, denoting by  $x$  the root of  $P$  in the interval, either  $Q(x) = 0$  or the sign of  $Q$  on the interval is everywhere equal to the sign of  $Q(x)$ .  $\square$

**Complexity analysis:**

Note first that the binary tree produced by Algorithm 10.6 is isolating for the polynomial  $PQ$ . Thus its number of nodes is at most  $O(p(\tau + \log_2(d)))$ , by Proposition 10.50.

The computation of  $G$  takes  $O(p^2)$  arithmetic operations, as well as the computation of  $b(G, \ell, m)$ ,  $b(\bar{P}, \ell, m)$ , and  $b(\bar{P}, m, r)$ .

We skip the details on the bit length as they are very similar to the ones in the binary complexity analysis of Algorithm 10.4 (Real Root Isolation).

Finally, the estimate for the binary complexity of the algorithm is  $O(p^5(\tau + \log_2(p))^2)$ .  $\square$

*Remark 10.57.* Note that it is easy to design variants to Algorithm 10.6, Algorithm 10.7, Algorithm 10.8 using Descartes’ isolation technique rather than Casteljau’s, with the same complexity.  $\square$

*Remark 10.58.* Similarly to Remark 10.52, using Remark 8.7, it is possible to compute the output of Algorithm 10.5 (Descartes' Real Root Isolation) with complexity  $\tilde{O}(p^2\tau)$  and binary complexity  $\tilde{O}(p^4\tau^2)$ . The same remark applies for the Descartes' variants of the other algorithms in this section: it is possible to compute the output of Algorithm 10.6 and Algorithm 10.7 with complexity  $\tilde{O}(p^2\tau)$  and binary complexity  $\tilde{O}(p^4\tau^2)$  and the output of Algorithm 10.8 with complexity  $\tilde{O}(s^2p^2\tau)$  and binary complexity  $\tilde{O}(s^2p^4\tau^2)$  using Descartes' variant (see Remark 10.57).  $\square$

We indicate now how to compare the roots of two polynomials in  $\mathbb{R}$ .

*Algorithm 10.7. [Comparison of Real Roots]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:** a polynomial  $P$  and a polynomial  $Q$  in  $\mathbb{Z}[X]$ .
- **Output:** a isolating list for the zeroes of  $\{P, Q\}$  in  $\mathbb{R}$  and  $\{P, Q\}$ .
- **Binary complexity:**  $O(p^5(\tau + \log_2(p))^2)$ , where  $p$  is a bound on the degree of  $P$ , and  $\tau$  is a bound on the bitsize of the coefficients of  $P$ .
- **Procedure:**
  - Compute  $\ell$  such that  $(-2^N, 2^N)$  contains the roots of  $P$  and  $Q$  using Lemma 10.2.
  - Isolate the roots of  $P$  (resp.  $Q$ ) using Algorithm 10.4 and perform the sign determination for  $Q$  (resp.  $P$ ) at these roots using Algorithm 10.6. Merge these two lists by taking the point or the interval of smallest length in case of non-empty intersection.

**Proof of correctness:** The algorithm terminates since  $\mathbb{R}$  is archimedean. Its correctness follows from Theorem 10.37. Note that because of the dichotomy process, any two elements of  $L(P)$  and  $L(Q)$  are either disjoint or one is included in the other.  $\square$

**Complexity analysis:** Follows from the binary complexity analysis of Algorithm 10.6 (Sign at a Real Root).  $\square$

Finally, we are able, given a finite set of univariate polynomials, to describe the real roots of these polynomials as well as points in the intervals they define.

*Algorithm 10.8. [Real Univariate Sample Points]*

- **Structure:** the ring  $\mathbb{Z}$ .
- **Input:** a finite set of univariate polynomials  $\mathcal{P}$  with coefficients in  $\mathbb{Z}$ .
- **Output:** an isolating list with signs for the roots of  $\mathcal{P}$  in  $\mathbb{R}$  and  $\mathcal{P}$ , an element between each two consecutive roots of elements of  $\mathcal{P}$ , an element of  $\mathbb{R}$  smaller than all these roots, and an element of  $\mathbb{R}$  greater than all these roots.
- **Binary complexity:**  $O(s^2p^5(\tau + \log_2(p))^2)$ , where  $p$  is a bound on the degree of  $P$ , and  $\tau$  is a bound on the bitsize of the coefficients of  $P$ .

• **Procedure:**

- For every pair  $P, Q$  of elements of  $\mathcal{P}$  perform Algorithm 10.7.
- Compute a rational point in between two consecutive roots using the isolating sets.
- Compute a rational point smaller than all these roots and a rational point greater than all the roots of polynomials in  $\mathcal{P}$  using Lemma 10.2.

**Proof of correctness:** The algorithm terminates since  $\mathbb{R}$  is archimedean. Its correctness follows from Theorem 10.37. □

**Complexity analysis:** Follows from the binary complexity analysis of Algorithm 10.7 (Comparison of Real Roots). □

### 10.3 Sign Determination

We consider now a general real closed field  $\mathbb{R}$ , not necessarily archimedean. Note that the approximation of the elements of  $\mathbb{R}$  by rational numbers cannot be performed anymore. Our aim is to give a method for determining the sign conditions realized by a family of polynomials on a finite set  $Z$  of points in  $\mathbb{R}^k$ .

This general method will be applied in two special cases: the zero set of a univariate polynomial in  $\mathbb{R}$  in this chapter and the zero set of a zero-dimensional polynomial system in  $\mathbb{R}^k$  later in the book.

Let  $Z$  be a finite subset of  $\mathbb{R}^k$ . We denote

$$\begin{aligned} \text{Reali}(P = 0, Z) &= \{x \in Z \mid P(x) = 0\}, \\ \text{Reali}(P > 0, Z) &= \{x \in Z \mid P(x) > 0\}, \\ \text{Reali}(P < 0, Z) &= \{x \in Z \mid P(x) < 0\}, \end{aligned}$$

and  $c(P = 0, Z)$ ,  $c(P > 0, Z)$ ,  $c(P < 0, Z)$  the corresponding numbers of elements. The Tarski-query of  $P$  for  $Z$  is

$$\text{TaQ}(P, Z) = \sum_{x \in Z} \text{sign}(Q(x)) = c(P > 0, Z) - c(P < 0, Z).$$

We consider the computation of  $\text{TaQ}(P, Z)$  as a basic black box. We have already seen several algorithms for computing it when  $Q \in \mathbb{R}[X]$ , and  $Z = \text{Zer}(Q, \mathbb{R})$  (Algorithms 9.2 and 9.5). Later in the book, we shall see other algorithms for the multivariate case.

Consider  $\mathcal{P} = P_1, \dots, P_s$ , a finite list of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ .

Let  $\sigma$  be a sign condition on  $\mathcal{P}$ , i.e. an element of  $\{0, 1, -1\}^{\mathcal{P}}$ . The **realization of the sign condition  $\sigma$  on  $Z$**  is

$$\text{Reali}(\sigma, Z) = \{x \in Z \mid \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma(P)\}.$$

The cardinality of  $\text{Reali}(\sigma, Z)$  is denoted  $c(\sigma, Z)$ .

We write  $\text{SIGN}(\mathcal{P}, Z)$  for the list of sign conditions realized by  $\mathcal{P}$  on  $Z$ , i.e. the list of  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$  such that  $\text{Reali}(\sigma, Z)$  is non-empty, and  $c(\mathcal{P}, Z)$  for the corresponding list of cardinals  $c(\sigma, Z) = \#(\text{Reali}(\sigma, Z))$  for  $\sigma \in \text{SIGN}(\mathcal{P}, Z)$ .

Our aim is to determine  $\text{SIGN}(\mathcal{P}, Z)$ , and, more precisely, to compute the numbers  $c(\mathcal{P}, Z)$ . The only information we are going to use to compute  $\text{SIGN}(\mathcal{P}, Z)$  is the Tarski-query of products of elements of  $\mathcal{P}$ .

A method for sign determination in the univariate case was already presented in Chapter 2 (Section 2.3). This method can be immediately generalized to the multivariate case, as we see now.

Given  $\alpha \in \{0, 1, 2\}^{\mathcal{P}}$ , we write  $\sigma^\alpha$  for  $\prod_{P \in \mathcal{P}} \sigma(P)^{\alpha(P)}$ , and  $\mathcal{P}^\alpha$  for  $\prod_{P \in \mathcal{P}} P^{\alpha(P)}$ , with  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$ . When  $\text{Reali}(\sigma, Z) \neq \emptyset$ , the sign of  $\mathcal{P}^\alpha$  is fixed on  $\text{Reali}(\sigma, Z)$  and is equal to  $\sigma^\alpha$  with the understanding that  $0^0 = 1$ .

We order the elements of  $\mathcal{P}$  so that  $\mathcal{P} = \{P_1, \dots, P_s\}$ . As in Chapter 2, we order  $\{0, 1, 2\}^{\mathcal{P}}$  lexicographically. We also order  $\{0, 1, -1\}^{\mathcal{P}}$  lexicographically (with  $0 < 1 < -1$ ).

Given  $A = \alpha_1, \dots, \alpha_m$ , a list of elements of  $\{0, 1, 2\}^{\mathcal{P}}$  with  $\alpha_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_m$ , we define

$$\begin{aligned} \mathcal{P}^A &= \mathcal{P}^{\alpha_1}, \dots, \mathcal{P}^{\alpha_m}, \\ \text{TaQ}(\mathcal{P}^A, Z) &= \text{TaQ}(\mathcal{P}^{\alpha_1}, Z), \dots, \text{TaQ}(\mathcal{P}^{\alpha_m}, Z). \end{aligned}$$

Given  $\Sigma = \sigma_1, \dots, \sigma_n$ , a list of elements of  $\{0, 1, -1\}^{\mathcal{P}}$ , with  $\sigma_1 <_{\text{lex}} \dots <_{\text{lex}} \sigma_n$ , we define

$$\begin{aligned} \text{Reali}(\Sigma, Z) &= \text{Reali}(\sigma_1, Z), \dots, \text{Reali}(\sigma_n, Z), \\ c(\Sigma, Z) &= c(\sigma_1, Z), \dots, c(\sigma_n, Z). \end{aligned}$$

The **matrix of signs of  $\mathcal{P}^A$  on  $\Sigma$**  is the  $m \times n$  matrix  $\text{Mat}(A, \Sigma)$  whose  $i, j$ -th entry is  $\sigma_j^{\alpha_i}$ .

**Proposition 10.59.** *If  $\cup_{\sigma \in \Sigma} \text{Reali}(\sigma, Z) = Z$ , then*

$$\text{Mat}(A, \Sigma) \cdot c(\Sigma, Z) = \text{TaQ}(\mathcal{P}^A, Z).$$

**Proof:** This is obvious since the  $(i, j)$ -th entry of  $\text{Mat}(\mathcal{P}^A, \Sigma)$  is  $\sigma_j^{\alpha_i}$ . □

When the matrix  $\text{Mat}(A, \Sigma)$  is invertible, we can compute  $c(\Sigma, Z)$  from  $\text{TaQ}(\mathcal{P}^A, Z)$ .

Note also that when  $\mathcal{P} = \{P\}$ ,  $A = \{0, 1, 2\}^{\{P\}}$ , and  $\Sigma = \{0, 1, -1\}^{\{P\}}$ , the conclusion of Proposition 10.59 is

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P=0, Z) \\ c(P>0, Z) \\ c(P<0, Z) \end{bmatrix} = \begin{bmatrix} \text{TaQ}(1, Z) \\ \text{TaQ}(P, Z) \\ \text{TaQ}(P^2, Z) \end{bmatrix}. \tag{10.13}$$

This is a generalization to  $Z$  of Equation (2.6) which had been stated for the set of zeroes of a univariate polynomial.

In order to compute each  $c(\sigma, Z)$  knowing all  $\text{TaQ}(\mathcal{P}^\alpha, Z)$ , we take  $A = \{0, 1, 2\}^{\mathcal{P}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{P}}$ .

As in Chapter 2, Notation 2.71 (Total matrix of signs), let  $M_s$  be the  $3^s \times 3^s$  matrix defined inductively by

$$M_1 = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix}$$

and

$$M_{t+1} = M_t \otimes M_1.$$

We generalize Proposition 2.72 and obtain

**Proposition 10.60.** *Let  $\mathcal{P}$  be a set of polynomials with  $s$  elements,  $A = \{0, 1, 2\}^{\mathcal{P}}$ , and  $\Sigma = \{0, 1, -1\}^{\mathcal{P}}$  ordered lexicographically. Then,*

$$\text{Mat}(A, \Sigma) = M_s.$$

**Proof:** The proof is by induction on  $s$ . If  $s=1$ , the claim is Equation (10.13). If the claim holds for  $s$ , it holds also for  $s+1$  given the definitions of  $M_{s+1}$  and  $\text{Mat}(\mathcal{P}^A, \Sigma)$ , and the orderings on  $A = \{0, 1, 2\}^{\mathcal{P}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{P}}$ .  $\square$

As a consequence:

**Corollary 10.61.**

$$M_s \cdot c(\Sigma, Z) = \text{TaQ}(\mathcal{P}^A, Z).$$

The preceding results give the following algorithm for sign determination, by using repeatedly the Tarski-query black box.

*Algorithm 10.9. [Naive Sign Determination]*

- **Input:** a finite subset  $Z \subset \mathbb{R}^k$  with  $r$  elements and a finite list  $\mathcal{P} = P_1, \dots, P_s$  of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ .
- **Output:** the list of sign conditions realized by  $\mathcal{P}$  on  $Z$ ,  $\text{SIGN}(\mathcal{P}, Z)$ .
- **Blackbox:** For a polynomial  $P$ , the Tarski-query  $\text{TaQ}(P, Z)$ .
- **Complexity:**  $3^s$  calls to the Tarski-query black box.
- **Procedure:**
  - Define  $A = \{0, 1, 2\}^{\mathcal{P}}$  and  $\Sigma = \{0, 1, -1\}^{\mathcal{P}}$ , ordered lexicographically.
  - Call the Tarski-query black box  $3^s$  times with input the elements of  $\mathcal{P}^A$  to obtain  $\text{TaQ}(\mathcal{P}^A, Z)$ . Solve the  $3^s \times 3^s$  system

$$M_s \cdot c(\Sigma, Z) = \text{TaQ}(\mathcal{P}^A, Z)$$

to obtain the vector  $c(\Sigma, Z)$  of length  $3^s$ . Output the set of sign conditions  $\sigma$  with  $c(\sigma, Z) \neq 0$ .

**Complexity analysis:** The number of calls to the Tarski-query black box is  $3^s$ . The calls to the Tarski-query black box are done for polynomials which are products of at most  $s$  polynomials of the form  $P$  or  $P^2$ ,  $P \in \mathcal{P}$ .  $\square$

To avoid the exponential number of calls to the Tarski-query black box in Algorithm 10.9 (Naive Sign Determination), notice that  $\#(\text{SIGN}(\mathcal{P}, Z)) \leq \#(Z)$ , so that the number of realizable sign conditions does not exceed  $\#(Z)$ . We are now going to determine the non-empty sign conditions inductively getting rid of the empty sign conditions at each step of the computation, in order to control the size of the data we manipulate.

Let  $Z \subset \mathbb{R}^k$  be a finite set, and let  $\mathcal{P}$  be a finite list of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ . A list  $A$  of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  is **adapted to sign determination for  $\mathcal{P}$  on  $Z$**  if the matrix of signs of  $\mathcal{P}^A$  over  $\text{SIGN}(\mathcal{P}, Z)$  is invertible.

*Example 10.62.* Consider the set of polynomials  $\{P\}$ . In this case,  $\{0, 1, 2\}^{\{P\}}$  can be identified with  $\{0, 1, 2\}$ . Note that when  $Z$  is non-empty,  $\text{SIGN}(\{P\}, Z)$  is also non-empty.

- If  $\text{SIGN}(\{P\}, Z) = \{0, 1, -1\}$ ,  $0, 1, 2$  is adapted to sign determination for  $\{P\}$  on  $Z$ , since  $\{P\}^{0,1,2} = 1, P, P^2$ , and the matrix of signs of  $1, P, P^2$  over  $0, 1, -1$  is

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix},$$

which is invertible.

- If  $\text{SIGN}(\{P\}, Z) = \{1, -1\}$  (resp.  $\{0, 1\}$ , resp.  $\{0, -1\}$ ),  $0, 1$  is adapted to sign determination for  $\{P\}$  on  $Z$ , since  $\{P\}^{0,1} = 1, P$  and the matrix of signs of  $1, P$  over  $1, -1$  (resp.  $0, 1$ , resp.  $0, -1$ ) is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad (\text{resp. } \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \text{ resp. } \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix}),$$

which is invertible.

- If  $\text{SIGN}(\{P\}, Z) = \{0\}$  (resp.  $\{1\}$ , resp.  $\{-1\}$ ),  $0$  is adapted to sign determination for  $\{P_i\}$  on  $Z$ , since  $\{P\}^0 = 1$  and the matrix of signs of  $1$  over  $0$  (resp.  $1$ , resp.  $-1$ ) is  $\begin{bmatrix} 1 \end{bmatrix}$ , which is invertible.

$\square$

Let  $Z \subset \mathbb{R}^k$  be a finite set,  $\mathcal{P}$  be a finite list of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ . We now describe a method for determining a list of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  adapted to sign determination for  $\mathcal{P}$  on  $Z$  from the set  $\text{SIGN}(\mathcal{P}, Z)$ . The definition of this list  $\text{Ada}(\mathcal{P}, Z)$  is by induction on the number of elements of  $\mathcal{P}$ .

Before defining  $\text{Ada}(\mathcal{P}, Z)$  we need the following definition.

**Definition 10.63. [Extension of a sign condition]** A sign condition  $\tau \in \text{SIGN}(\{P\} \cup \mathcal{P}, Z)$  extends  $\sigma \in \text{SIGN}(\mathcal{P}, Z)$  if  $\sigma(Q) = \tau(Q), Q \in \mathcal{P}$ .  $\square$

We now define  $\text{Ada}(\mathcal{P}, Z)$ .

**Definition 10.64. [Adapted family]**

- If  $\mathcal{P} = \{P\}$ ,
  - if  $\#(\text{SIGN}(\{P\}, Z)) = 3$ , define  $\text{Ada}(\{P\}, Z) = 0, 1, 2$ ,
  - if  $\#(\text{SIGN}(\{P\}, Z)) = 2$ , define  $\text{Ada}(\{P\}, Z) = 0, 1$ ,
  - if  $\#(\text{SIGN}(\{P\}, Z)) = 1$ , define  $\text{Ada}(\{P\}, Z) = 0$ .
- If  $\mathcal{P} = \{P\} \cup \mathcal{Q}$ , let  $\text{SIGN}(\mathcal{Q}, Z)_2$  be the subset of  $\text{SIGN}(\mathcal{Q}, Z)$  of sign conditions  $\sigma$  such that there are at least two distinct sign conditions of  $\text{SIGN}(\mathcal{P}, Z)$  extending  $\sigma$ , and  $\text{SIGN}(\mathcal{Q}, Z)_3$  be the subset of  $\text{SIGN}(\mathcal{Q}, Z)$  of sign conditions  $\sigma$  such that there are three distinct sign conditions of  $\text{SIGN}(\mathcal{P}, Z)$  extending  $\sigma$ . Let

$$Z_2 = \bigcup_{\sigma \in \text{SIGN}(\mathcal{Q}, Z)_2} \text{Reali}(\sigma, Z),$$

$$Z_3 = \bigcup_{\sigma \in \text{SIGN}(\mathcal{Q}, Z)_3} \text{Reali}(\sigma, Z).$$

Note that

$$\begin{aligned} \text{SIGN}(\mathcal{Q}, Z_2) &= \text{SIGN}(\mathcal{Q}, Z)_2, \\ \text{SIGN}(\mathcal{Q}, Z_3) &= \text{SIGN}(\mathcal{Q}, Z)_3. \end{aligned}$$

For  $\alpha \in \{0, 1, 2\}$  and  $\beta \in \{0, 1, 2\}^{\mathcal{Q}}$ , we define  $\alpha \times \beta \in \{0, 1, 2\}^{\mathcal{P}}$  by

$$\begin{cases} (\alpha \times \beta)(P) = \alpha(P), \\ (\alpha \times \beta)(Q) = \beta(Q) \text{ if } Q \in \mathcal{Q}. \end{cases}$$

Define

$$\text{Ada}(\mathcal{P}, Z) = 0 \times \text{Ada}(\mathcal{Q}, Z), 1 \times \text{Ada}(\mathcal{Q}, Z_2), 2 \times \text{Ada}(\mathcal{P}, Z_3).$$

$\square$

**Proposition 10.65.** *The list  $\text{Ada}(\mathcal{P}, Z)$  is adapted to sign determination for  $\mathcal{P}$  on  $Z$ .*

**Proof:** The proof proceeds by induction on the number of elements of  $\mathcal{P}$ . The claim is true for  $\mathcal{P} = \{P\}$ , as seen in Example 10.62.

If  $\mathcal{P} = \{P\} \cup \mathcal{Q}$ , we want to prove that

$$\text{Mat}(\text{Ada}(\mathcal{P}, Z), \text{SIGN}(\mathcal{P}, Z))$$

is invertible. Denoting by  $C_\tau$  its column indexed by  $\tau$ , consider a zero linear combination of its columns:

$$\sum_{\tau \in \text{SIGN}(\mathcal{P}, Z)} \lambda_\tau C_\tau = 0.$$



We want to prove that all  $\lambda_\tau$  are zero.

If  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_3$ , we denote by  $\sigma_1 <_{\text{lex}} \sigma_2 <_{\text{lex}} \sigma_3$  the sign conditions of  $\text{SIGN}(\mathcal{P}, Z)$  extending  $\sigma$ .

Similarly, if  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_2 \setminus \text{SIGN}(\mathcal{Q}, Z)_3$ , we denote by

$$\sigma_1 <_{\text{lex}} \sigma_2$$

the sign conditions of  $\text{SIGN}(\mathcal{P}, Z)$  extending  $\sigma$ .

Finally if  $\sigma \in \text{SIGN}(\mathcal{Q}, Z) \setminus \text{SIGN}(\mathcal{Q}, Z)_2$ , we denote by  $\sigma_1$  the sign condition of  $\text{SIGN}(\mathcal{P}, Z)$  extending  $\sigma$ .

Since by induction hypothesis  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{Q}, Z))$  is invertible,

$$\begin{aligned} \lambda_{\sigma_1} &= 0, & \text{for every } \sigma \in \text{SIGN}(\mathcal{Q}, Z) \setminus \text{SIGN}(\mathcal{Q}, Z)_2, \\ \lambda_{\sigma_1} + \lambda_{\sigma_2} &= 0, & \text{for every } \sigma \in \text{SIGN}(\mathcal{Q}, Z)_2 \setminus \text{SIGN}(\mathcal{Q}, Z)_3, \\ \lambda_{\sigma_1} + \lambda_{\sigma_2} + \lambda_{\sigma_3} &= 0, & \text{for every } \sigma \in \text{SIGN}(\mathcal{Q}, Z)_3. \end{aligned}$$

By induction hypothesis, the matrix  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{Q}, Z_2))$  is invertible, then  $\sigma_1(P) \lambda_{\sigma_1} - \sigma_2(P) \lambda_{\sigma_2} = 0$ , for every  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_2 \setminus \text{SIGN}(\mathcal{Q}, Z)_3$  and  $\lambda_{\sigma_2} - \lambda_{\sigma_3} = 0$ , for every  $\text{SIGN}(\mathcal{Q}, Z)_3$ . Thus  $\lambda_{\sigma_1} = \lambda_{\sigma_2} = 0$ , for every  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_2 \setminus \text{SIGN}(\mathcal{Q}, Z)_3$ . Finally, using again the induction hypothesis,  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z_3), \text{SIGN}(\mathcal{Q}, Z_3))$  is invertible, then  $\lambda_{\sigma_2} + \lambda_{\sigma_3} = 0$  for every  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_3$ . Thus  $\lambda_{\sigma_1} = \lambda_{\sigma_2} = \lambda_{\sigma_3} = 0$  for every  $\sigma \in \text{SIGN}(\mathcal{Q}, Z)_3$ .

This proves that the matrix

$$\text{Mat}(\text{Ada}(\mathcal{P}, Z), \text{SIGN}(\mathcal{P}, Z))$$

is invertible. □

**Lemma 10.66.** *Let  $Z' \subset Z$ ,  $r = \#(\text{SIGN}(\mathcal{P}, Z))$ ,  $r' = \#(\text{SIGN}(\mathcal{P}, Z'))$ . The matrix  $\text{Mat}(\text{Ada}(\mathcal{P}, Z'), \text{SIGN}(\mathcal{P}, Z'))$  coincides with the matrix obtained by extracting from  $\text{Mat}(\text{Ada}(\mathcal{P}, Z), \text{SIGN}(\mathcal{P}, Z'))$  its  $r'$  first linearly independent rows.*

**Proof:** The proof is by induction on the number of elements of  $\mathcal{P}$ .

The claim is clearly true is  $\mathcal{P} = \{P\}$ .

Suppose now that  $\mathcal{P} = \{P\} \cup \mathcal{Q}$  and the claim holds for  $\mathcal{Q}$ .

Note that by Definition 10.64,  $\text{Ada}(\mathcal{P}, Z')$  is a sublist of  $\text{Ada}(\mathcal{P}, Z)$ , so that the rank of  $\text{Mat}(\text{Ada}(\mathcal{P}, Z), \text{SIGN}(\mathcal{P}, Z'))$  is  $r'$ .

Similarly,

$$\begin{aligned} r'_1 &= \#(\text{SIGN}(\mathcal{Q}, Z')) = \text{Rank}(\text{Mat}(\text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{Q}, Z'))), \\ r'_2 &= \#(\text{SIGN}(\mathcal{Q}, Z')_2) = \text{Rank}(\text{Mat}(\text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{Q}, Z'))), \\ r'_3 &= \#(\text{SIGN}(\mathcal{Q}, Z')_3) = \text{Rank}(\text{Mat}(\text{Ada}(\mathcal{Q}, Z_3), \text{SIGN}(\mathcal{Q}, Z'))). \end{aligned}$$

It follows immediately that,

$$\begin{aligned} \text{Rank}(\text{Mat}(0 \times \text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{P}, Z'))) &\leq r'_1, \\ \text{Rank}(\text{Mat}(1 \times \text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{P}, Z'))) &\leq r'_2, \\ \text{Rank}(\text{Mat}(2 \times \text{Ada}(\mathcal{Q}, Z_3), \text{SIGN}(\mathcal{P}, Z'))) &\leq r'_3. \end{aligned}$$

Finally, since  $r'_1 + r'_2 + r'_3 = r$ ,

$$\begin{aligned} \text{Rank}(\text{Mat}(0 \times \text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{P}, Z'))) &= r'_1, \\ \text{Rank}(\text{Mat}(1 \times \text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{P}, Z'))) &= r'_2, \\ \text{Rank}(\text{Mat}(2 \times \text{Ada}(\mathcal{Q}, Z_3), \text{SIGN}(\mathcal{P}, Z'))) &= r'_3, \end{aligned}$$

and the first  $r'$  linearly independent rows of  $\text{Mat}(\text{Ada}(\mathcal{P}, Z), \text{SIGN}(\mathcal{P}, Z'))$  consist of  $r'_1$  rows of  $\text{Mat}(0 \times \text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{P}, Z'))$ ,  $r'_2$  linearly independent rows of  $\text{Mat}(1 \times \text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{P}, Z'))$  and  $r'_3$  linearly independent rows of  $\text{Mat}(2 \times \text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{P}, Z'))$ . The corresponding  $r'_1$  (resp.  $r'_2$ , resp.  $r'_3$ ) rows of  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{Q}, Z'))$  (resp.  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z_2), \text{SIGN}(\mathcal{Q}, Z'_2))$ , resp.  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z_3), \text{SIGN}(\mathcal{Q}, Z'_3))$ ) are linearly independent and are the rows indexed by  $\text{Ada}(\mathcal{Q}, Z')$  (resp.  $\text{Ada}(\mathcal{Q}, Z'_2)$ , resp.  $\text{Ada}(\mathcal{Q}, Z'_3)$ ) by the induction hypothesis. The claim follows from Definition 10.64. □

*Algorithm 10.10.* [Adapted Family]

- **Input:** the set  $\text{SIGN}(\{P\} \cup \mathcal{Q}, Z)$ , the list  $\text{Ada}(\mathcal{Q}, Z)$ .
- **Output:** the list  $\text{Ada}(\{P\} \cup \mathcal{Q}, Z)$ .
- **Procedure:**
  - If  $\mathcal{Q} = \emptyset$ ,
    - if  $\#(\text{SIGN}(\{P\}, Z)) = 3$ , define  $\text{Ada}(\{P\}, Z) = 0, 1, 2$ ,
    - if  $\#(\text{SIGN}(\{P\}, Z)) = 2$ , define  $\text{Ada}(\{P\}, Z) = 0, 1$ ,
    - if  $\#(\text{SIGN}(\{P\}, Z)) = 1$ , define  $\text{Ada}(\{P\}, Z) = 0$ .
  - Using the notation in Definition 10.64, let

$$\begin{aligned} r_1 &= \#(\text{SIGN}(\mathcal{Q}, Z)), \\ r_2 &= \#(\text{SIGN}(\mathcal{Q}, Z)_2), \\ r_3 &= \#(\text{SIGN}(\mathcal{Q}, Z)_3). \end{aligned}$$

Then  $\#(\text{SIGN}(\{P\} \cup \mathcal{Q}, Z)) = r_1 + r_2 + r_3$ .

Consider the matrix  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{Q}, Z)_2)$  and extract from it the first  $r_2$  linearly independent rows, which correspond to a sublist  $A_2$  of  $\text{Ada}(\mathcal{Q}, Z)$ .

Similarly, consider the matrix  $\text{Mat}(\text{Ada}(\mathcal{Q}, Z), \text{SIGN}(\mathcal{Q}, Z)_3)$  and extract from it the first  $r_3$  linearly independent rows which correspond to a sublist  $A_3$  of  $\text{Ada}(\mathcal{Q}, Z)$ .

Output

$$\text{Ada}(\{P\} \cup \mathcal{Q}, Z) = 0 \times \text{Ada}(\mathcal{Q}, Z), 1 \times A_2, 2 \times A_3.$$

**Correctness of Algorithm 10.10 :** Follows immediately from Definition 10.64 and Lemma 10.66. □

**Notation 10.67. [Sign determination]** Let  $\mathcal{P} = P_1, \dots, P_s$ . For  $1 \leq i \leq s$  we define  $\mathcal{P}_i = P_i, \dots, P_s$ . For  $\sigma \in \{0, 1, -1\}$  and  $\tau \in \{0, 1, -1\}^{\mathcal{P}^{i+1}}$ , we define  $\sigma \wedge \tau \in \{0, 1, -1\}^{\mathcal{P}^i}$  by

$$\begin{cases} (\sigma \wedge \tau)(P_i) = \sigma(P_i) \\ (\sigma \wedge \tau)(P) = \tau(P) \quad \text{if } P \in \mathcal{P}_{i+1}, \end{cases}$$

If  $\Sigma = \sigma_1, \dots, \sigma_m$  is a list of elements of  $\{0, 1, -1\}$  with  $\sigma_1 <_{\text{lex}} \dots <_{\text{lex}} \sigma_m$  and  $T = \tau_1, \dots, \tau_n$  is a list of element of  $\{0, 1, -1\}^{\mathcal{P}^i}$  with  $\tau_1 <_{\text{lex}} \dots <_{\text{lex}} \tau_n$ , then  $\Sigma \wedge T$  is the list

$$\sigma_1 \wedge \tau_1 <_{\text{lex}} \dots <_{\text{lex}} \sigma_1 \wedge \tau_n <_{\text{lex}} \dots <_{\text{lex}} \sigma_m \wedge \tau_1 <_{\text{lex}} \dots <_{\text{lex}} \sigma_m \wedge \tau_n.$$

For  $\alpha \in \{0, 1, 2\}$  and  $\beta \in \{0, 1, 2\}^{\mathcal{P}^{i+1}}$ , we define  $\alpha \times \beta \in \{0, 1, 2\}^{\mathcal{P}^i}$  by

$$\begin{cases} (\alpha \times \beta)(P_i) = \alpha, \\ (\alpha \times \beta)(P) = \beta(P) \quad \text{if } P \in \mathcal{P}_{i+1}. \end{cases}$$

If  $A = \alpha_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_m$  and  $B = \beta_1 <_{\text{lex}} \dots <_{\text{lex}} \beta_n$  are lists of elements of  $\{0, 1, 2\}$  and  $\{0, 1, 2\}^{\mathcal{P}^{i+1}}$  we define  $A \times B$  to be the list

$$\alpha_1 \times \beta_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_1 \times \beta_n <_{\text{lex}} \dots <_{\text{lex}} \alpha_m \times \beta_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_m \times \beta_n$$

in  $\{0, 1, 2\}^{\mathcal{P}^i}$ .

The list  $\mathcal{P}_i^{A \times B}$  is defined to be

$$P_i^{\alpha_1} \mathcal{P}_{i+1}^{\beta_1}, \dots, P_i^{\alpha_1} \mathcal{P}_{i+1}^{\beta_n}, \dots, P_i^{\alpha_m} \mathcal{P}_{i+1}^{\beta_1}, \dots, P_i^{\alpha_m} \mathcal{P}_{i+1}^{\beta_n}. \quad \square$$

Recall that the matrix of signs of  $\mathcal{P}^B$  (of length  $m$ ) on  $\Sigma$  (of length  $n$ ) is the  $m \times n$  matrix  $\text{Mat}(B, \Sigma)$  whose  $i, j$ -th entry is  $\sigma_j^{\alpha_i}$ , and that  $\text{TaQ}(\mathcal{P}^B, Z)$  is the vector  $\text{TaQ}(\mathcal{P}^{\beta_1}, Z), \dots, \text{TaQ}(\mathcal{P}^{\beta_m}, Z)$ . Using Notation 2.69 (Tensor product) we have

**Proposition 10.68.** *If  $\cup_{\sigma \in \Sigma} \text{Reali}(\sigma, Z) = Z$   $A = 0, 1, 2$ , and  $T = \{0, 1, -1\}$ , let*

$$(\text{Mat}(A, T) \otimes \text{Mat}(B, \Sigma)) \cdot c(T \wedge \Sigma, Z) = \text{TaQ}(\mathcal{P}_i^{A \times B}, Z).$$

**Proof:** Immediate from Proposition 10.59. □

We are now ready for the Sign Determination algorithm.

*Algorithm 10.11. [Sign Determination]*

- **Input:** a finite subset  $Z \subset \mathbb{R}^k$  with  $r$  elements and a finite list  $\mathcal{P} = P_1, \dots, P_s$  of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ .
- **Output:** the list of sign conditions realized by  $\mathcal{P}$  on  $Z$ ,  $\text{SIGN}(\mathcal{P}, Z)$ .
- **Blackbox:** for a polynomial  $P$ , the Tarski-query  $\text{TaQ}(P, Z)$ .
- **Complexity:**  $1 + 2sr$  calls to the to the Tarski-query black box.

• **Procedure:**

- Compute  $r = \text{TaQ}(1, Z)$  using the Tarski-query black box with input 1. If  $r = 0$ , output  $\emptyset$ .
- Let  $\mathcal{P}_i = P_i, \dots, P_s$ . We are going to determine iteratively, for  $i$  from  $s$  to 1,  $\text{SIGN}(\mathcal{P}_i, Z)$  the non-empty sign conditions for  $\mathcal{P}_i$  on  $Z$ . More precisely, we are going to compute  $\text{SIGN}(\mathcal{P}_i, Z)$  and  $\text{Ada}(\mathcal{P}_i, Z)$ , starting from  $\text{SIGN}(\mathcal{P}_{i+1}, Z)$  and  $\text{Ada}(\mathcal{P}_{i+1}, Z)$ .
- For  $i$  from  $s$  to 1,
  - Determine  $\text{SIGN}(P_i, Z)$ , the list of sign conditions realized by  $P_i$  on  $Z$ , and a list  $B_i$  of elements in  $\{0, 1, 2\}$  adapted to sign determination for  $P_i$  on  $Z$  as follows:
    - Use the Tarski-query black box with inputs  $P_i$  and  $P_i^2$  to determine  $\text{TaQ}(P_i, Z)$  and  $\text{TaQ}(P_i^2, Z)$ .
    - From these values, using the equality

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} c(P_i = 0, Z) \\ c(P_i > 0, Z) \\ c(P_i < 0, Z) \end{bmatrix} = \begin{bmatrix} \text{TaQ}(1, Z) \\ \text{TaQ}(P_i, Z) \\ \text{TaQ}(P_i^2, Z) \end{bmatrix},$$

compute  $c(P_i = 0, Z)$ ,  $c(P_i > 0, Z)$  and  $c(P_i < 0, Z)$  and output  $\text{SIGN}(P_i, Z)$ .

- If  $r(P_i) = \#(\text{SIGN}(P_i, Z)) = 3$ , let  $B_i = 0, 1, 2$ .
- If  $r(P_i) = \#(\text{SIGN}(P_i, Z)) = 2$ , let  $B_i = 0, 1$ .
- If  $r(P_i) = \#(\text{SIGN}(P_i, Z)) = 1$ , let  $B_i = 0$ .
- Define  $M_i = \text{Mat}(B_i, \text{SIGN}(P_i, Z))$ .
- If  $i = s$ , define  $\text{SIGN}(\mathcal{P}_s, Z) := \text{SIGN}(P_s, Z)$ ,  $\text{Ada}(\mathcal{P}_s, Z) := B_s$ .
- If  $i < s$ , Compute  $\text{SIGN}(\mathcal{P}_i, Z)$ , the list of sign conditions realized by  $\mathcal{P}_i$  on  $Z$ , as follows:
  - Use the Tarski-query black box with input the elements of  $\mathcal{P}_i^{B_i \times \text{Ada}(\mathcal{P}_{i+1}, Z)}$  to determine  $d' = \text{TaQ}(\mathcal{P}_i^{B_i \times \text{Ada}(\mathcal{P}_{i+1}, Z)}, Z)$ .
  - Take the matrix

$$M'_i := \text{Mat}(\text{Ada}(\mathcal{P}_{i+1}, Z), \text{SIGN}(\mathcal{P}_{i+1}, Z)) \otimes M_i.$$

Compute the list  $c' = c(\text{SIGN}(P_i, Z) \wedge \text{SIGN}(\mathcal{P}_{i+1}, Z))$  from the equality  $M'_i \cdot c' = d'$  by inverting  $M'_i$ . Compute  $\text{SIGN}(\mathcal{P}_i, Z)$ , removing from  $\text{SIGN}(P_i, Z) \wedge \text{SIGN}(\mathcal{P}_{i+1}, Z)$  the sign conditions with empty realization, which correspond to the zeroes in  $c'$ .

- Call Algorithm 10.10 (Adapted family) with input  $\text{SIGN}(\mathcal{P}_i, Z)$  and  $\text{Ada}(\mathcal{P}_{i+1}, Z)$ , and compute  $\text{Ada}(\mathcal{P}_i, Z)$ .
- Output  $\text{SIGN}(\mathcal{P}, Z) = \text{SIGN}(\mathcal{P}_1, Z)$ .

*Remark 10.69.* We denote by  $B(\text{SIGN}(\mathcal{P}, Z)) \subset \{0, 1, 2\}^{\mathcal{P}}$  the set constructed inductively as follows:

$$\begin{aligned} B(\text{SIGN}(\mathcal{P}_s, Z)) &= \{0, 1, 2\}_1 \\ B(\text{SIGN}(\mathcal{P}_i, Z)) &= B(\text{SIGN}(\mathcal{P}_{i+1}, Z)) \cup \{0, 1, 2\}_i \cup B_i \times \text{Ada}(\mathcal{P}_{i+1}, Z), \end{aligned}$$

denoting by  $\{0, 1, 2\}_i$  the subset of  $\{0, 1, 2\}^{\mathcal{P}}$  with three elements defined by

$$\alpha \in \{0, 1, 2\}_i \Leftrightarrow \alpha(j) = 0 \forall j \neq i,$$

and identifying  $\alpha \in \{0, 1, 2\}^{\mathcal{P}_i}$  with  $\alpha' \in \{0, 1, 2\}^{\mathcal{P}}$  such that

$$\alpha'(P_j) = \alpha(P_j), j \geq i, \alpha'(P_j) = 0, j < i,$$

using the notation of Algorithm 10.11 (Sign Determination). It is easy to see that  $B(\text{SIGN}(\mathcal{P}, Z))$  is nothing but the list of elements  $\alpha \in \{0, 1, 2\}^{\mathcal{P}}$  such that the Tarski-query of  $P^\alpha$  has been computed in Algorithm 10.11 (Sign Determination). Using Algorithm 10.10 (Adapted family), it is clear that  $B(\text{SIGN}(\mathcal{P}, Z))$  can be determined from  $\text{SIGN}(\mathcal{P}, Z)$ .  $\square$

**Proof of correctness of Algorithm 10.11:** It follows from Corollary 10.68 and the correctness of Algorithm 10.10 (Adapted family).  $\square$

Before discussing the complexity of the Sign Determination Algorithm, we first give an example.

*Example 10.70.* Consider

$$\begin{aligned} P &= (X^3 - 1)(X^2 - 9), \\ Z &= \text{Zer}(P, R), \\ P_1 &= X - 2, \\ P_2 &= X + 1, \\ P_3 &= X. \end{aligned}$$

The call to the Tarski-query black box with input 1 determines  $\text{TaQ}(1, Z) = 3$ . So  $P$  has 3 real roots (which is not a real surprise).

The call to the Tarski-query black box with inputs  $P_3$  and  $P_3^2$  determines  $\text{TaQ}(P_3, Z) = 1$  and  $\text{TaQ}(P_3^2, Z) = 3$ . Thus

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P_3 = 0, Z) \\ c(P_3 > 0, Z) \\ c(P_3 < 0, Z) \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 3 \end{bmatrix},$$

which means, after solving the system, that  $P$  has

$$\begin{cases} 0 & \text{root with } P_3 = 0 \\ 2 & \text{roots with } P_3 > 0. \\ 1 & \text{root with } P_3 < 0 \end{cases}$$

Hence  $c(P_3 = 0, Z) = 0$ . So we have  $\text{SIGN}(P_3, Z) = 1, -1$  and

$$\text{Ada}(\mathcal{P}_3, Z) = B_3 = 0, 1.$$

The matrix  $\text{Mat}(\text{Ada}(\mathcal{P}_3, Z), \text{SIGN}(\mathcal{P}_3, Z))$  of signs of 1,  $P_3$  on 1,  $-1$  is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

We now consider  $\mathcal{P}_2 = P_2, P_3$ .

The call to the Tarski-query black box with inputs  $P_2$  and  $P_2^2$  determines  $\text{TaQ}(P_2, Z) = 1, \text{TaQ}(P_2^2, Z) = 3$ . Hence,

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P_2=0, Z) \\ c(P_2>0, Z) \\ c(P_2<0, Z) \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 3 \end{bmatrix},$$

which means, after solving the system, that  $P$  has

$$\begin{cases} 0 & \text{root with } P_2 = 0 \\ 2 & \text{roots with } P_2 > 0 \\ 1 & \text{root with } P_2 < 0 \end{cases}.$$

Hence  $c(P_2 = 0, Z) = 0$ . So we have  $\text{SIGN}(P_2, Z) = 1, -1$  and  $B_2 = 0, 1$ . The matrix  $M_2$  of signs of  $1, P_2$  on the sign conditions  $1, -1$  is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

The call to the Tarski-query black box with input  $P_2 P_3$  yields  $\text{TaQ}(P_2 P_3, Z)$ , which is equal to 3. Hence we have

$$M'_2 = \text{Mat}(\text{Ada}(\mathcal{P}_3, Z), \text{SIGN}(\mathcal{P}_3, Z)) \otimes M_2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix},$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P_2 > 0 \wedge P_3 > 0, Z) \\ c(P_2 > 0 \wedge P_3 < 0, Z) \\ c(P_2 < 0 \wedge P_3 > 0, Z) \\ c(P_2 < 0 \wedge P_3 < 0, Z) \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 1 \\ 3 \end{bmatrix}.$$

Solving the system we find that  $P$  has

$$\begin{cases} 2 & \text{roots with } P_2 > 0 \text{ and } P_3 > 0 \\ 0 & \text{roots with } P_2 > 0 \text{ and } P_3 < 0 \\ 0 & \text{roots with } P_2 < 0 \text{ and } P_3 > 0 \\ 1 & \text{root with } P_2 < 0 \text{ and } P_3 < 0 \end{cases}.$$

So we have  $\text{SIGN}(\mathcal{P}_2, Z) = (1, 1), (-1, -1)$ . There is no sign condition on  $P_3$  which is partitioned by sign conditions on  $P_2$ , so  $\text{Ada}(\mathcal{P}_2, Z) = (0, 0), (1, 0)$ . The matrix  $\text{Mat}(\text{Ada}(\mathcal{P}_2, Z), \text{SIGN}(\mathcal{P}_2, Z))$  of signs of  $1, P_3$  on the sign conditions  $(1, 1), (-1, -1)$  is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Finally we consider  $\mathcal{P} = P_1, P_2, P_3$ .

The call to the Tarski-query black box with inputs  $P_1$  and  $P_1^2$  determines  $\text{TaQ}(P_1, Z) = -1, \text{TaQ}(P_1^2, Z) = 3$ . Hence  $c(P_1 = 0, Z) = 0$ . So,

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P_1 = 0, Z) \\ c(P_1 > 0, Z) \\ c(P_1 < 0, Z) \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \\ 3 \end{bmatrix},$$

which means, after solving the system, that  $P$  has

$$\begin{cases} 0 & \text{root with } P_1 = 0 \\ 1 & \text{root with } P_1 > 0 \\ 2 & \text{roots with } P_1 < 0 \end{cases}.$$

So we have  $\text{SIGN}(P_1, Z) = \{1, -1\}, B_1 = \{0, 1\}$ . The matrix  $M_1$  of signs of 1,  $P_1$  on 1,  $-1$  is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

The call to the Tarski-query black box with input  $P_1 P_3$  yields  $\text{TaQ}(P_1 P_3, Z)$  which is equal to 1. Hence we have

$$M'_1 = \text{Mat}(\text{Ada}(\mathcal{P}_2, Z), \text{SIGN}(\mathcal{P}_2, Z)) \otimes M_1 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix},$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(P_1 > 0, P_2 > 0, P_3 > 0, Z) \\ c(P_1 > 0, P_2 > 0, P_3 < 0, Z) \\ c(P_1 < 0, P_2 < 0, P_3 > 0, Z) \\ c(P_1 < 0, P_2 < 0, P_3 < 0, Z) \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \\ 1 \\ 1 \end{bmatrix}.$$

Solving the system, we find that  $P$  has

$$\begin{cases} 1 & \text{root with } P_1 > 0 \text{ and } P_2 > 0 \text{ and } P_3 > 0 \\ 0 & \text{root with } P_1 > 0 \text{ and } P_2 < 0 \text{ and } P_3 < 0 \\ 1 & \text{root with } P_1 < 0 \text{ and } P_2 > 0 \text{ and } P_3 > 0 \\ 1 & \text{root with } P_1 < 0 \text{ and } P_2 < 0 \text{ and } P_3 < 0 \end{cases}.$$

So we have  $\text{SIGN}(\mathcal{P}) = \{(1, 1, 1), (-1, 1, 1), (-1, -1, -1)\}$ . □

In order to study the complexity of the Algorithm 10.11 (Sign Determination) we need the following proposition.

**Proposition 10.71.** *Let  $Z$  be a finite subset of  $\mathbb{R}^k$  and  $r = \#(Z)$ . For every  $\alpha \in \text{Ada}(\mathcal{P}, Z)$ , the number  $\#(\{P \in \mathcal{P} \mid \alpha(P) \neq 0\})$  is at most  $\log_2(r)$ .*

We need the following definition. Let  $\alpha$  and  $\beta$  be elements of  $\{0, 1, 2\}^{\mathcal{P}}$ . We say that  $\beta$  **precedes**  $\alpha$  if for every  $P \in \mathcal{P}$ ,  $\beta(P) \neq 0$  implies  $\beta(P) = \alpha(P)$ . Note that if  $\beta$  precedes  $\alpha$ , then  $\beta <_{\text{lex}} \alpha$ .

The proof of Proposition 10.71 is based on the following lemma.

**Lemma 10.72.** *If  $\beta$  precedes  $\alpha$  and  $\alpha \in \text{Ada}(\mathcal{P}, Z)$  then  $\beta \in \text{Ada}(\mathcal{P}, Z)$ .*

**Proof:** The proof is by induction on the number of elements of  $\mathcal{P}$ . The claim is obvious for  $\mathcal{P} = \{P\}$ .

If  $\mathcal{P} = \{P\} \cup \mathcal{Q}$ ,  $\alpha \in \{0, 1, 2\}^{\mathcal{P}}$  we denote by  $\alpha'$  the element of  $\{0, 1, 2\}^{\mathcal{Q}}$  such that  $\alpha'(P) = \alpha(P)$ ,  $P \in \mathcal{Q}$ . Note that, by definition of  $\text{Ada}(\mathcal{P}, Z)$ , if  $\alpha' \notin \text{Ada}(\mathcal{Q}, Z)$ ,  $\alpha \notin \text{Ada}(\mathcal{P}, Z)$ .

Suppose that  $\beta$  precedes  $\alpha$  and that  $\beta \notin \text{Ada}(\mathcal{P}, Z)$ . There are several cases to consider:

- If  $\alpha(P) = \beta(P) = 1$ ,  $\beta' \notin \text{Ada}(\mathcal{Q}, Z_2)$  by Definition 10.64. By induction hypothesis,  $\alpha' \notin \text{Ada}(\mathcal{Q}, Z_2)$  and  $\alpha = 1 \times \alpha' \notin \text{Ada}(\mathcal{P}, Z)$  again by Definition 10.64.
- If  $\alpha(P) = \beta(P) = 2$ ,  $\beta' \notin \text{Ada}(\mathcal{Q}, Z_3)$  by Definition 10.64. By induction hypothesis,  $\alpha' \notin \text{Ada}(\mathcal{Q}, Z_3)$  and  $\alpha = 2 \times \alpha' \notin \text{Ada}(\mathcal{P}, Z)$  again by Definition 10.64.
- If  $\beta(P) = 0$ ,  $\beta' \notin \text{Ada}(\mathcal{Q}, Z)$  by Definition 10.64, thus  $\alpha' \notin \text{Ada}(\mathcal{Q}, Z)$  by induction hypothesis, and  $\alpha \notin \text{Ada}(\mathcal{P}, Z)$  by Definition 10.64. □

**Proof of Proposition 10.71:**

Let  $\alpha$  be such that  $\#\{P \in \mathcal{P} \mid \alpha(P) \neq 0\} = k$ . Since the number of elements  $\beta$  of  $\{0, 1, 2\}^{\mathcal{P}}$  preceding  $\alpha$  is  $2^k$ , and the total number of polynomials in  $A_s$  is at most  $r$ , we have  $2^k \leq r$  and  $k \leq \log_2(r)$ . So, the claim follows immediately from the Lemma 10.72. □

**Complexity analysis:** There are  $s$  steps in Algorithm 10.11 (Sign Determination). In each step, the number of calls to the Tarski-query black box is bounded by  $2r$ . Indeed, in Step  $i$ , there are at most  $3r_{i-1}$  Tarski-queries to compute and  $r_{i-1}$  of these Tarski-queries have been determined in Step  $i - 1$ . So, in Step  $i$ , there are at most  $2r_{i-1}$  Tarski-queries to determine. The total number of calls to the the Tarski-query black box is bounded by  $1 + 2sr$ . The calls to the Tarski-query black box are done for polynomials which are product of at most  $\log_2(r)$  products of polynomials of the form  $P$  or  $P^2$ ,  $P \in \mathcal{P}$  by Proposition 10.71. □

Note that we did not count the complexity of performing the linear algebra involved in the algorithm. This is because when we consider particular ways of realization the Tarski-query black box later, we bound only the number of arithmetic operations in the ring. Since the complexity of linear algebra is polynomial in the size of the matrix, the maximum size of the matrices is  $3r$ , and their entries are 0, 1 or  $-1$ , taking into account the linear algebra part of the algorithm would not change the linearity in  $s$  and the polynomial time in  $r$  character of the algorithm.



We finally describe how to get a family adapted to sign determination from  $\text{SIGN}(\mathcal{P}, Z)$  using Algorithm 10.10 (Adapted family).

*Algorithm 10.12.* **[Family adapted to Sign Determination]**

- **Input:** the set  $\text{SIGN}(\mathcal{P}, Z)$ .
- **Output:** a list  $\text{Ada}(\mathcal{P}, Z)$  of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  adapted to sign determination for  $\mathcal{P}$  on  $Z$ .
- **Procedure:**  
 Let  $\mathcal{P} = P_1, \dots, P_s$ ,  $\mathcal{P}_i = P_i, \dots, P_s$ ,  $i = 1, \dots, s$ . Note that  $\text{SIGN}(\mathcal{P}_i, Z)$  can be obtained from  $\text{SIGN}(\mathcal{P}, Z)$  by forgetting the  $i - 1$  first signs of elements of  $\text{SIGN}(\mathcal{P}, Z)$ .
  - If  $\#(\text{SIGN}(\mathcal{P}_s, Z)) = 3$ , define  $\text{Ada}(\mathcal{P}_s, Z) = 0, 1, 2$ .
  - If  $\#(\text{SIGN}(\mathcal{P}_s, Z)) = 2$ , define  $\text{Ada}(\mathcal{P}_s, Z) = 0, 1$ .
  - If  $\#(\text{SIGN}(\mathcal{P}_s, Z)) = 1$ , define  $\text{Ada}(\mathcal{P}_s, Z) = 0$ .
  - For  $i$  from  $s - 1$  to 1, apply Algorithm 10.10 (Adapted family) to  $\text{SIGN}(\mathcal{P}_i, Z)$ ,  $\text{Ada}(\mathcal{P}_{i+1}, Z)$  to obtain  $\text{Ada}(\mathcal{P}_i, Z)$ .
  - Output  $\text{Ada}(\mathcal{P}_1, Z)$ .

We can now describe in a more specific way how the Tarski-query black box can be implemented in the univariate case.

*Algorithm 10.13.* **[Univariate Sign Determination]**

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $R$ .
- **Input:** a non-zero univariate polynomial  $Q$  and a list  $\mathcal{P}$  of univariate polynomials with coefficients in  $D$ . Let  $Z = \text{Zer}(Q, R)$ .
- **Output:** the list of sign conditions realized by  $\mathcal{P}$  on  $Z$ ,  $\text{SIGN}(\mathcal{P}, Z)$ , and a list  $A$  of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  adapted to sign determination for  $\mathcal{P}$  on  $Z$ .
- **Complexity:**  $O(sp^2(p + q \log_2(p)))$ , where  $s$  is a bound on the number of polynomials in  $\mathcal{P}$ ,  $p$  is a bound on the degree of  $Q$  and  $q$  is a bound on the degree of the polynomials in  $\mathcal{P}$ .
- **Procedure:** Perform Algorithm 10.11 (Sign Determination), using as Tarski-query black box Algorithm 9.5 (Univariate Tarski-query). Products of elements of  $\mathcal{Q}$  are reduced modulo  $P$  each time a multiplication is performed.

**Complexity analysis:** According to the complexity of Algorithm 10.11 (Sign Determination), the number of calls to the Tarski-query black box is bounded by  $1 + 2sp$ , since  $r \leq p$ . The calls to the Tarski-query black box are done for  $P$  and polynomials of degree at most  $q$ . The complexity is thus  $O(sp^2(p + q))$ , using the complexity of Algorithm 9.5 (Univariate Tarski-query).

When  $Q$  and  $P \in \mathcal{P}$  are in  $\mathbb{Z}[X]$  with coefficients of bitsize bounded by  $\tau$ , the bitsize of the integers in the operations performed by the algorithm are bounded by  $O((p + q \log_2(p))(\tau + \log_2(p + q \log_2(p))))$ , according to Proposition 8.44.  $\square$

*Remark 10.73.* Using Remark 9.2, it is possible to compute the output of Algorithm 10.13 (Univariate Sign Determination) with complexity  $\tilde{O}(sp(p+q))$  and binary complexity  $\tilde{O}(sp(p+q)^2\tau)$ .  $\square$

### 10.4 Roots in a Real Closed Field

We consider here too a general real closed field  $\mathbb{R}$ , not necessarily archimedean. In such a field, it is not possible to perform real root isolation and to approximate roots by rational numbers. In order to characterize and compute the roots of a polynomial, in a sense made precise in this section, we are going to use Proposition 2.28 (Thom encoding) and the preceding sign determination method.

Let  $P \in \mathbb{R}[X]$  and  $\sigma \in \{0, 1, -1\}^{\text{Der}(P)}$ , a sign condition on the set  $\text{Der}(P)$  of derivatives of  $P$ . By Definition 2.29, the sign condition  $\sigma$  is a Thom encoding of  $x \in \mathbb{R}$  if  $\sigma(P) = 0$  and  $\sigma$  is the sign condition taken by the set  $\text{Der}(P)$  at  $x$ . We say that  $x$  is specified by  $\sigma$ . Given a Thom encoding  $\sigma$ , we denote by  $x(\sigma)$  the root of  $P$  in  $\mathbb{R}$  specified by  $\sigma$ .

The **ordered list of Thom encodings of  $P$**  is the ordered list  $\sigma_1, \dots, \sigma_r$  of Thom encodings of the roots  $x(\sigma_1) < \dots < x(\sigma_r)$  of  $P$ .

The ordered list of Thom encodings of a univariate polynomial can be obtained using sign determination as follows.

*Algorithm 10.14.* **[Thom Encoding]**

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $\mathbb{R}$ .
- **Input:** a non-zero polynomial  $P \in D[X]$  of degree  $p$ .
- **Output:** the ordered list of Thom encodings of the roots of  $P$  in  $\mathbb{R}$ .
- **Complexity:**  $O(p^4 \log_2(p))$ .
- **Procedure:** Apply Algorithm 10.13 (Univariate Sign Determination) to  $P$  and its derivatives  $\text{Der}(P')$ . Order the Thom encodings using Proposition 2.28.

**Complexity analysis:** The complexity is  $O(p^4 \log_2(p))$  using the complexity of Algorithm 10.13 (Univariate Sign Determination), since Algorithm 10.13 is called with a family of at most  $p$  polynomials of degree at most  $p$ .

When  $P \in \mathbb{Z}[X]$ , with coefficients of bitsize bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O(p \log_2(p) (\tau + \log_2(p)))$  according to Proposition 8.44.  $\square$

*Remark 10.74.* When arithmetic operations are performed naively, it follows from the preceding complexity analysis, using Remark 8.4, that the binary complexity of Algorithm 10.14 (Thom Encodings) is thus

$$O(p^6 \log_2(p)^3 (\tau + \log_2(p))^2).$$

Note that from a binary complexity point of view, Algorithms 10.4 (Real Root Isolation) is preferable to Algorithm 10.14 (Thom Encodings). It turns out that, in practice as well, Algorithm 10.4 is much better, as the number of nodes in the isolation tree of Algorithm 10.4 (Real Root Isolation) is much smaller in most cases than its theoretical value  $O(p(\tau + \log_2(p)))$  given by Proposition 10.50. This is the reason why, even though it is less general than Algorithm 10.14 (Thom Encoding), Algorithm 10.4 (Real Root Isolation) is important.  $\square$

*Remark 10.75.* Using Remark 9.2, it is possible to compute the output of Algorithm 10.14 in complexity  $\tilde{O}(p^3)$  and binary complexity  $\tilde{O}(p^4 \tau)$ . Similarly the output of Algorithm 10.15 can be computed in complexity  $\tilde{O}(p^2(p+q))$  and binary complexity  $\tilde{O}(p^2(p+q)^2 \tau)$ , the output of Algorithm 10.16 and Algorithm 10.18 can be computed in complexity  $\tilde{O}(p^3)$  and binary complexity  $\tilde{O}(p^4 \tau)$ , and the output of Algorithm 10.17 and Algorithm 10.19 in complexity  $\tilde{O}(s^2 p^3)$  and binary complexity  $\tilde{O}(s^2 p^4 \tau)$ .  $\square$

*Remark 10.76.* The Thom Encoding algorithm is based on the Sign Determination algorithm which is in turn based on the Signed subresultant Algorithm. This algorithm uses exact divisions and is valid only in an integral domain, and not in a general ring. In a ring, the algorithm computing determinants indicated in Remark 8.19 can always be used for computing the signed subresultant coefficients, and hence the Thom encoding. The complexity obtained is  $p^{O(1)}$  arithmetic operations in the ring  $D$  of coefficients of  $P$ , which is sufficient for the complexity estimates proved in later chapters.  $\square$

*Algorithm 10.15.* **[Sign at the Roots in a Real Closed Field]**

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $R$ .
- **Input:** a polynomial  $P \in D[X]$  of degree  $p$  and a polynomial  $Q \in D[X]$  of degree  $q$ , the list  $\text{Thom}(P)$  of Thom encodings of the set  $Z$  of roots of  $P$  in  $R$ .
- **Output:** for every  $\sigma \in \text{Thom}(P)$  specifying the root  $x$  of  $P$ , the sign  $Q(x)$ .
- **Complexity:**  $O(p^2(p \log_2(p) + q))$ .
- **Procedure:**
  - Determine the non-empty sign conditions  $\text{SIGN}(Q, Z)$  for  $Q$  and the list  $\text{Ada}(Q)$  of elements in  $\{0, 1, 2\}$  adapted to sign determination using Algorithm 10.11 (Sign Determination).
  - Construct from the list  $\text{Thom}(P)$  of Thom encodings of the roots of  $P$  the list  $\text{Ada}(\text{Der}(P'))$  of elements in  $\{0, 1, 2\}^{\text{Der}(P')}$  adapted to sign determination using Algorithm (Family adapted for sign determination) 10.12
  - Determine the non-empty sign conditions for  $\text{Der}(P')$ ,  $Q$  as follows:
    - Compute the list of Tarski-queries

$$d' = \text{TaQ}((Q, \text{Der}(P'))^{\text{Ada}(Q) \times \text{Ada}(\text{Der}(P'))}, Z).$$

– Let  $M = \text{Mat}(\text{Der}(P'), \text{SIGN}(\text{Der}(P'), Z))$  and

$$M' = \text{Mat}(\text{Ada}(Q), \text{SIGN}(Q, Z)) \otimes M.$$

Compute the list  $c' = c(\text{SIGN}(Q, Z) \wedge \text{SIGN}(\text{Der}(P'), Z))$  from the equality

$$M' \cdot c' = d'$$

by inverting  $M'$ . Output using the non-zero entries of  $c'$  the signs of  $Q(x(\sigma))$ ,  $\sigma \in \text{SIGN}(\text{Der}(P'), Z)$ .

**Proof of correctness:** This is a consequence of Proposition 2.28 since the number of non-zero elements in  $c'$  is exactly  $r = c(\text{Der}(P'))$ . □

**Complexity analysis:** The complexity is  $O(p^2(p \log_2(p) + q))$  since there are at most  $3p$  calls to Algorithm 9.5 (Univariate Tarski-query) for polynomials of degree  $p$  and  $p \log_2(p) + q$ .

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , and the bitsizes of the coefficients of  $P$  and  $Q$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $(\tau + \log_2(p + q)) O(p \log_2(p + q) + q)$ , using the complexity analysis of Algorithm 9.5 (Univariate Tarski-query). □

It is also possible to compare the roots of two polynomials in a real closed field by a similar method.

Let  $\mathcal{P}$  be a finite subset of  $\mathbb{R}[X]$ . The **ordered list of Thom encodings of  $\mathcal{P}$**  is the ordered list  $\sigma_1, \dots, \sigma_r$  of Thom encoding of elements of

$$Z = \{x \in \mathbb{R} \mid \bigvee_{P \in \mathcal{P}} P(x) = 0\} = \{x(\sigma_1) < \dots < x(\sigma_r)\}.$$

*Algorithm 10.16.* **[Comparison of Roots in a Real Closed Field]**

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $\mathbb{R}$ .
- **Input:** two non-zero polynomials  $P$  and  $Q$  in  $D[X]$  of degree  $p$ .
- **Output:** the ordered list of the Thom encodings of  $\{P, Q\}$ .
- **Complexity:**  $O(p^4 \log_2(p))$ .
- **Procedure:** Apply Algorithm 10.13 (Univariate Sign Determination) to  $P$  and  $\text{Der}(P')$ ,  $\text{Der}(Q)$ , then to  $Q$  and  $\text{Der}(Q')$ ,  $\text{Der}(P)$ . Compare the roots using Proposition 2.28.

**Complexity analysis:** The complexity is  $O(p^4 \log_2(p))$  since we call Algorithm 10.13 (Univariate Sign determination) twice, each time with a family of at most  $2p$  polynomials of degree at most  $p$ .

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , with coefficients of bitsize bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O(p \log_2(p) (\tau + \log_2(p)))$  according to Proposition 8.44. □

Finally, we are able, given a finite set of univariate polynomials, to describe the ordered list of real roots of these polynomials.

*Algorithm 10.17. [Partition of a Line]*

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $R$ .
- **Input:** a finite family  $\mathcal{P} \subset D[X]$ .
- **Output:** the ordered list of the roots of  $\mathcal{P}$ , described by Thom encodings.
- **Complexity:**  $O(s^2 p^4 \log_2(p))$ , where  $p$  is a bound on the degree of the elements of  $\mathcal{P}$ , and  $s$  a bound on the number of elements of  $\mathcal{P}$ .
- **Procedure:** Characterize all the roots of the polynomials of  $\mathcal{P}$  in  $R$  using Algorithm 10.14 (Thom Encoding). Using Algorithm 10.16, compare these roots for every couple of polynomials in  $\mathcal{D}$ . Output the ordered list of Thom encodings of  $\mathcal{P}$ .

**Complexity analysis:** Since there are  $O(s^2)$  pairs of polynomials to consider, the complexity is clearly bounded by  $O(s^2 p^4 \log_2(p))$ , using the complexity of Algorithms 10.16.

When  $\mathcal{P} \subset \mathbb{Z}[X]$  and the coefficients of  $P \in \mathcal{P}$  are of bitsize bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O(p \log_2(p) (\tau + \log_2(p)))$  according to Proposition 8.44.  $\square$

It is also possible, using the same techniques, to find a point between two elements of  $R$  specified by Thom encodings.

*Algorithm 10.18. [Intermediate Points]*

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $R$ .
- **Input:** two non-zero univariate polynomials  $P$  and  $Q$  in  $R[X]$  of degree bounded by  $p$ .
- **Output:** Thom encodings specifying values  $y$  in intervals between two consecutive roots of  $P$  and  $Q$ .
- **Complexity:**  $O(p^4 \log_2(p))$ .
- **Procedure:** Compute the Thom encodings of the roots of  $(PQ)'$  in  $R$  using Algorithm 10.14 (Thom Encoding) and compare them to the roots of  $P$  and  $Q$  using Algorithm 10.16. Keep one intermediate point between two consecutive roots of  $PQ$ .

**Proof of correctness:** Let  $y$  be a root of  $P$  and  $z$  be a root of  $Q$ . Then there is a root of  $(PQ)'$  in  $(y, z)$  by Rolle's theorem (Proposition 2.22).  $\square$

**Complexity analysis:** The complexity is clearly bounded by  $O(p^4 \log_2(p))$  using the complexity analysis of Algorithms 10.14 and 10.16.

When  $P$  and  $Q$  are in  $\mathbb{Z}[X]$ , with coefficients of bitsize bounded by  $\tau$ , the bitsize of the integers in the operations performed by the algorithm are bounded by  $O(p \log_2(p) (\tau + \log_2(p)))$  according to Proposition 8.44.  $\square$

*Remark 10.77.* Note that Algorithm 10.18 (Intermediate Points) can also be used to produce intermediate points between zeros of one polynomial by setting  $Q = 1$ .  $\square$

Finally we are able, given a finite set of univariate polynomials, to describe the real roots of these polynomials as well as points between consecutive roots.

Given a family  $\mathcal{P}$  of univariate polynomials, an **ordered list of sample points for  $\mathcal{P}$**  is an ordered list  $L$  of Thom encodings  $\sigma$  specifying the roots of the polynomials of  $\mathcal{P}$  in  $\mathbb{R}$ , an element between two such consecutive roots, an element of  $\mathbb{R}$  smaller than all these roots, and an element of  $\mathbb{R}$  greater than all these roots. Moreover  $\sigma$ , appears before  $\tau$  in  $L$  if and only if  $x(\sigma) \leq x(\tau)$ . The sign of  $Q(x(\sigma))$  is also output for every  $Q \in \mathcal{P}, \sigma \in L$ .

*Algorithm 10.19. [Univariate Sample Points]*

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite subset  $\mathcal{P} \subset D[X]$ .
- **Output:** an ordered list of sample points for  $\mathcal{P}$ .
- **Complexity:**  $O(s^2 p^4 \log_2(p))$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $p$  is a bound on the degree of the elements of  $\mathcal{P}$ .
- **Procedure:** Characterize all the roots of the polynomials in  $\mathbb{R}$  using Algorithm 10.14 (Thom Encoding). Using Algorithm 10.16, compare these roots for every couple of polynomials in  $\mathcal{P}$ . Compute a itemize of a point in each interval between the roots by Algorithm 10.18 (Intermediate Points). Order all these Thom encodings and keep only one intermediate point in each open interval between roots of polynomials in  $\mathcal{P}$ . Use Proposition 10.1 to find a polynomial of degree 1 with coefficients in  $D$  whose root is smaller (resp. larger) than any root of any polynomial in  $\mathcal{P}$ .

**Complexity analysis:** Since there are  $O(s^2)$  pairs of polynomials to consider, the complexity is clearly bounded by  $O(s^2 p^4 \log_2(p))$ , using the complexity of Algorithms 10.16 and 10.18.

When  $\mathcal{P} \subset \mathbb{Z}[X]$  and the coefficients of  $P \in \mathcal{P}$  are of bitsize bounded by  $\tau$ , the bitsizes of the integers in the operations performed by the algorithm are bounded by  $O(p \log_2(p) (\tau + \log_2(p)))$  according to Proposition 8.44.  $\square$

## 10.5 Bibliographical Notes

The real root isolation method goes back to Vincent [162] and has been studied by Uspensky [158]. Bernstein's polynomials are important in Computer Aided Design [57], they have been used in real root isolation in [101]. The algorithm for computing the Bernstein coefficients described in this chapter was discovered by De Casteljaou, an engineer. The complexity analysis of the real root isolation algorithm is due to [55].

The basic idea of the sign determination algorithm appears in [23]. The use of Thom encodings for characterizing real roots appears in [50].

---

## Cylindrical Decomposition Algorithm

The cylindrical decomposition method described in Chapter 5 can be turned into algorithms for solving several important problems.

The first problem is the **general decision problem** for the theory of the reals. The general decision problem is to design a procedure to decide the truth or falsity of a sentence  $\Phi$  of the form  $(\text{Qu}_1 X_1) \dots (\text{Qu}_k X_k) F(X_1, \dots, X_k)$ , where  $\text{Qu}_i \in \{\exists, \forall\}$  and  $F(X_1, \dots, X_k)$  is a quantifier free formula.

The second problem is the **quantifier elimination problem**. We are given a formula  $\Phi(Y)$  of the form  $(\text{Qu}_1 X_1) \dots (\text{Qu}_k X_k) F(Y_1, \dots, Y_\ell, X_1, \dots, X_k)$ , where  $\text{Qu}_i \in \{\exists, \forall\}$  and  $F(Y_1, \dots, Y_\ell, X_1, \dots, X_k)$  is a quantifier free formula. The quantifier elimination problem is to output a quantifier free formula,  $\Psi(Y)$ , such that for any  $y \in \mathbb{R}^\ell$ ,  $\Phi(y)$  is true if and only if  $\Psi(y)$  is true.

The general decision problem is a special case of the quantifier elimination problem, corresponding to  $\ell = 0$ .

In Chapter 2, we have already proved that every formula is equivalent to a quantifier free formula (Theorem 2.77). The method used in the proof can in fact be turned into an algorithm, but if we performed the complexity analysis of this algorithm, we would get a tower of exponents of height linear in the number of variables. We decided not to develop the complexity analysis of the method for quantifier elimination presented in Chapter 2, since the algorithms described in this chapter and in Chapter 14 have a much better complexity.

In Section 11.1, we describe the Cylindrical Decomposition Algorithm. The degrees of the polynomials output by this algorithm are doubly exponential in the number of variables. A general outline is included in the first part of the section, and technical details on the lifting phase are included in the second part. In Section 11.2 we use the Cylindrical Decomposition Algorithm to decide the truth of a sentence. In Section 11.3, a variant of the Cylindrical Decomposition Algorithm makes it possible to perform quantifier elimination. In Section 11.4, we prove that the complexity of quantifier elimination is intrinsically doubly exponential. In Section 11.5, another variant of the Cylindrical Decomposition Algorithm is used to compute a stratification. In the two variable case, the Cylindrical Decomposition Algorithm is particularly

simple and is used for computing the topology of a real algebraic plane curve in Section 11.6. Finally in Section 11.7, a variant called Restricted Elimination is used to replace infinitesimal quantities by sufficiently small numbers.

## 11.1 Computing the Cylindrical Decomposition

### 11.1.1 Outline of the Method

We use the results in Section 5.1 of Chapter 5, in particular the definition of a cylindrical decomposition adapted to  $\mathcal{P}$  (see Definitions 5.1 and 5.5) and the properties of the set  $\text{Elim}_{X_k}(\mathcal{P})$  (see Notation 5.15).

We denote, for  $i = k - 1, \dots, 1$ ,

$$C_i(\mathcal{P}) = \text{Elim}_{X_{i+1}}(C_{i+1}(\mathcal{P})),$$

with  $C_k(\mathcal{P}) = \mathcal{P}$ , so that  $C_i(\mathcal{P}) \subset \mathbb{R}[X_1, \dots, X_i]$ . The family

$$C(\mathcal{P}) = \bigcup_{i \leq k} C_i(\mathcal{P})$$

is the **cylindrifying family of polynomials associated to  $\mathcal{P}$** . It follows from the proof of Theorem 5.6 that the semi-algebraically connected components of the sign conditions on  $C(\mathcal{P})$  are the cells of a cylindrical decomposition adapted to  $\mathcal{P}$ .

The Cylindrical Decomposition Algorithm consists of two phases: in the first phase the cylindrifying family of polynomials associated to  $\mathcal{P}$  is computed and in the second phase the cells defined by these polynomials are used to define inductively, starting from  $i = 1$ , the cylindrical decomposition.

The computation of the cylindrifying family of polynomials associated to  $\mathcal{P}$  is based on the following Elimination Algorithm, computing the family of polynomials  $\text{Elim}_{X_k}(\mathcal{P})$  defined in 5.15, using Notation 1.16.

#### *Algorithm 11.1.* [Elimination]

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite list of variables  $X_1, \dots, X_k$ , a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , and a variable  $X_k$ .
- **Output:** a finite set  $\text{Elim}_{X_k}(\mathcal{P}) \subset D[X_1, \dots, X_{k-1}]$ . The set  $\text{Elim}_{X_k}(\mathcal{P})$  is such that the degree of  $P \in \mathcal{P}$  with respect to  $X_k$ , the number of real roots of  $P \in \mathcal{P}$ , and the number of real roots common to  $P \in \mathcal{P}$  and  $Q \in \mathcal{P}$  is fixed on every semi-algebraically connected component of the realization of each sign condition on  $\text{Elim}_{X_k}(\mathcal{P})$ .
- **Complexity:**  $s^2 d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$ , and  $d$  is a bound on the degrees of the elements of  $\mathcal{P}$ .



- **Procedure:** Place in  $\text{Elim}_{X_k}(\mathcal{P})$  the following polynomials, computed by Algorithm 8.21 (Signed Subresultant), using Remark 8.50, when they are not in  $D$ :
  - For  $P \in \mathcal{P}$ ,  $\deg_{X_k}(P) = p \geq 2$ ,  $R \in \text{Tru}(P)$ ,  $j = 0, \dots, \deg_{X_k}(R) - 2$ ,  $\text{sRes}_j(R, \partial R / \partial X_k)$ .
  - For  $R \in \text{Tru}(\mathcal{P})$ ,  $S \in \text{Tru}(\mathcal{P})$ ,
    - if  $\deg_{X_k}(R) > \deg_{X_k}(S)$ ,  $\text{sRes}_j(R, S)$ ,  $j = 0, \dots, \deg_{X_k}(S) - 1$ ,
    - if  $\deg_{X_k}(R) < \deg_{X_k}(S)$ ,  $\text{sRes}_j(S, R)$ ,  $j = 0, \dots, \deg_{X_k}(R) - 1$ ,
    - if  $\deg_{X_k}(R) = \deg_{X_k}(S)$ ,  $\text{sRes}_j(S, \bar{R})$ , with  $\bar{R} = \text{lcof}(S)R - \text{lcof}(R)S$ ,  $j = 0, \dots, \deg_{X_k}(\bar{R}) - 1$ .
  - For  $R \in \text{Tru}(\mathcal{P})$ ,  $\text{lcof}(R)$ .

**Proof of correctness:** The correctness follows from Theorem 5.16. □

**Complexity analysis of Algorithm 11.1:** Consider

$$D[X_1, \dots, X_k] = D[X_1, \dots, X_{k-1}][X_k].$$

There are  $O(s^2 d^2)$  subresultant sequences to compute, since there are  $O(s^2)$  couples of polynomials in  $\mathcal{P}$  and  $O(d)$  truncations for each polynomial to consider. Each of these subresultant sequence takes  $O(d^2)$  arithmetic operations in the integral domain  $D[X_1, \dots, X_{k-1}]$  according to the complexity analysis of Algorithm 8.21 (Signed Subresultant). The complexity is thus  $O(s^2 d^4)$  in the integral domain  $D[X_1, \dots, X_{k-1}]$ . There are  $O(s^2 d^3)$  polynomials output.

The degree with respect to  $X_1, \dots, X_{k-1}$  of the polynomials throughout these computations is bounded by  $2 d^2$  by Proposition 8.45. Since each multiplication and exact division of polynomials of degree  $2 d^2$  in  $k - 1$  variables costs  $O(d)^{4(k-1)}$  (see Algorithms 8.5 and 8.6), the final complexity is  $s^2 O(d)^{4k} = s^2 d^{O(k)}$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and output are bounded by  $\tau d^{O(k)}$ , using Proposition 8.46. □

*Example 11.1.* a) Let  $P = X_1^2 + X_2^2 + X_3^2 - 1$ . The output of Algorithm 11.1 (Elimination) with input the variable  $X_3$  and the set  $\mathcal{P} = \{P\}$  is (getting rid of irrelevant constant factors) the polynomial  $\text{sRes}_0(P, \partial P / \partial X_3) = X_1^2 + X_2^2 - 1$  (see Example 5.17).

b) Consider the two polynomials

$$P = X_2^2 - X_1(X_1 + 1)(X_1 - 2), Q = X_2^2 - (X_1 + 2)(X_1 - 1)(X_1 - 3).$$

The output of Algorithm 11.1 (Elimination) with input the variable  $Y$  and  $\mathcal{P} = \{P, Q\}$  contains three polynomials: the discriminant of  $P$  with respect to  $X_2$ ,

$$\text{sRes}_0(P, \partial P / X_2) = 4 X_1(X_1 + 1)(X_1 - 2),$$

the discriminant of  $Q$  with respect to  $Y$ ,

$$\text{sRes}_0(Q, \partial Q / \partial X_2) = 4(X_1 + 2)(X_1 - 1)(X_1 - 3),$$

and the resultant of  $P$  and  $Q$  with respect to  $Y$ ,

$$\text{sRes}_0(P, Q) = (-X_1^2 - 3X_1 + 6)^2,$$

since  $\text{sRes}_1(P, Q) = 0$  is a constant.  $\square$

Now we are ready to describe the two phases of the cylindrical decomposition method.

Let  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_k$  be a cylindrical decomposition of  $\mathbb{R}^k$ . A **cylindrical set of sample points** of  $\mathcal{S}$ ,  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_k$ , is a list of  $k$  sets such that

- for every  $i$ ,  $1 \leq i \leq k$ ,  $\mathcal{A}_i$  is a finite subset of  $\mathbb{R}^i$  which intersects every  $S \in \mathcal{S}_i$ ,
- for every  $i$ ,  $1 \leq i \leq k - 1$ ,  $\pi_i(\mathcal{A}_{i+1}) = \mathcal{A}_i$ , where  $\pi_i$  is the projection from  $\mathbb{R}^{i+1}$  to  $\mathbb{R}^i$  forgetting the last coordinate.

*Algorithm 11.2. [Cylindrical Decomposition]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite ordered list of variables  $X_1, \dots, X_k$ , and a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** a cylindrical set of sample points of a cylindrical decomposition  $\mathcal{S}$  adapted to  $\mathcal{P}$  and the sign of the elements of  $\mathcal{P}$  on each cell of  $\mathcal{S}_k$ .
- **Complexity:**  $(sd)^{O(1)k-1}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$ , and  $d$  is a bound on the degrees of the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Initialize  $C_k(\mathcal{P}) := \mathcal{P}$ .
  - Elimination phase: Compute  $C_i(\mathcal{P}) = \text{Elim}_{X_{i+1}}(C_{i+1}(\mathcal{P}))$ , for  $i = k - 1, \dots, 1$ , applying repeatedly  $\text{Elim}_{X_{i+1}}$  using Algorithm 11.1 (Elimination).
  - Lifting phase:
    - Compute the sample points of the cells in  $\mathcal{S}_1$  by characterizing the roots of  $C_1(\mathcal{P})$  and choosing a point in each interval they determine.
    - For every  $i = 2, \dots, k$ , compute the sample points of the cells of  $\mathcal{S}_i$  from the sample points of the cells in  $\mathcal{S}_{i-1}$  as follows: Consider, for every sample point  $x$  of a cell in  $\mathcal{S}_{i-1}$ , the list  $L$  of non-zero polynomials  $P_i(x, X_i)$  with  $P_i \in C_i(\mathcal{P})$ . Characterize the roots of  $L$  and choose a point in each interval they determine.
  - Output the sample points of the cells and the sign of  $P \in \mathcal{P}$  on the corresponding cells of  $\mathbb{R}^k$ .

We need to be more specific about how we describe and compute sample points. This will be explained fully in the next subsection.

**Proof of correctness:** The correctness of Algorithm 11.2 (Cylindrical Decomposition) follows from the proof of Theorem 5.6.  $\square$

**Complexity analysis of the Elimination phase.** Using the complexity analysis of Algorithm 11.1 (Elimination), if the input polynomials have degree  $D$ , the degree of the output is  $2(D^2)$  after one application of Algorithm 11.1 (Elimination). Thus, the degrees of the polynomials output after  $k - 1$  applications of Algorithm 11.1 (Elimination) are bounded by  $f(d, k - 1)$ , where  $f$  satisfies the recurrence relation

$$f(d, i) = 2f(d, i - 1)^2, f(d, 0) = d. \quad (11.1)$$

Solving the recurrence we get that  $f(d, k) = 2^{1+2+\dots+2^{k-2}}d^{2^{k-1}}$ , and hence the degrees of the polynomials in the intermediate computations and the output are bounded by  $2^{1+2+\dots+2^{k-2}}d^{2^{k-1}} = O(d)^{2^{k-1}}$ , which is polynomial in  $d$  and doubly exponential in  $k$ . A similar analysis shows that the number of polynomials output is bounded by  $(sd)^{3^{k-1}}$ , which is polynomial in  $s$  and  $d$  and doubly exponential in  $k$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the the intermediate computations and the output are bounded by  $\tau d^{O(1)k-1}$ , using the complexity analysis of Algorithm 11.1 (Elimination), which is performed  $k - 1$  times.  $\square$

*Example 11.2.* Let  $P = X_1^2 + X_2^2 + X_3^2 - 1$ . Continuing Example 5.17, we describe the output of the Cylindrical Decomposition Algorithm applied to  $\mathcal{P} = \{P\}$ .

We have

$$\begin{aligned} C_3(\mathcal{P}) &= \{X_1^2 + X_2^2 + X_3^2 - 1\}, \\ C_2(\mathcal{P}) &= \{X_1^2 + X_2^2 - 1\}, \\ C_1(\mathcal{P}) &= \{X_1^2 - 1\}. \end{aligned}$$

The sample points of  $\mathbb{R}$  consists of five points, corresponding to the two roots of  $X^2 - 1$  and one point in each of the three intervals they define: these are the semi-algebraically connected components of the realization of sign conditions defined by  $C_1(\mathcal{P})$ . We choose a sample point in each cell and obtain

$$\{(S_1, -2), (S_2, -1), (S_3, 0), (S_4, 1), (S_5, 2)\}.$$

The cells in  $\mathbb{R}^2$  are obtained by taking the semi-algebraically connected components of the realization of sign conditions defined by  $C_1(\mathcal{P}) \cup C_2(\mathcal{P})$ . There are thirteen such cells, listed in Example 5.4. The sample points in  $\mathbb{R}^2$  consist of thirteen points, one in each cell. The projection of a sample point in a cell of  $\mathbb{R}^2$  on its first coordinate is a point in a cell of  $\mathbb{R}$ . We choose a sample point in each cell and obtain

$$\begin{aligned} &\{(S_{1,1}, (-2, 0)), \\ &\quad (S_{2,1}, (-1, -1)), (S_{2,2}, (-1, 0)), (S_{2,3}, (-1, 1)), \\ &\quad (S_{3,1}, (0, -2)), (S_{3,2}, (0, -1)), (S_{3,3}, (0, 0)), (S_{3,4}, (0, 1)), (S_{3,5}, (0, 2)), \\ &\quad (S_{4,1}, (1, -1)), (S_{4,2}, (1, 0)), (S_{4,3}, (1, 1)), \\ &\quad (S_{5,1}, (2, 0))\}. \end{aligned}$$

The cells in  $\mathbb{R}^3$  are obtained by taking the semi-algebraically connected components of the realization of sign conditions defined by

$$C_1(\mathcal{P}) \cup C_2(\mathcal{P}) \cup C_3(\mathcal{P}).$$

There are twenty five such cells, listed in Example 5.4. The sample points in  $\mathbb{R}^3$  consist of twenty five points, one in each cell. The projection of a sample point in a cell of  $\mathbb{R}^3$  is a point in a cell of  $\mathbb{R}^2$ . We choose the following sample points and obtain, indicating the cell, its sample point and the sign of  $P$  at this sample point:

$$\begin{aligned} & \{(S_{1,1,1}, (-2, 0, 0), 1), \\ & \quad (S_{2,1,1}, (-1, -1, 0), 1), \\ (S_{2,2,1}, (-1, 0, -1), 1), & (S_{2,2,2}, (-1, 0, 0), 0), (S_{2,2,3}, (-1, 0, 1), 1), \\ & \quad (S_{2,3,1}, (-1, 1, 0), 1), \\ & \quad (S_{3,1,1}, (0, -2, 0), 1), \\ (S_{3,2,1}, (0, -1, -1), 1), & (S_{3,2,2}, (0, -1, 0), 0), (S_{3,2,3}, (0, -1, 1), 1), \\ & \quad (S_{3,3,1}, (0, 0, -2), 1), (S_{3,3,2}, (0, 0, -1), 0), \\ & \quad (S_{3,3,3}, (0, 0, 0), -1), \quad \square \\ & \quad (S_{3,3,4}, (0, 0, 1), 0), (S_{3,3,5}, (0, 0, 2), 1), \\ (S_{3,4,1}, (0, 1, -1), 1), & (S_{3,4,2}, (0, 1, 0), 0), (S_{3,4,3}, (0, 1, 1), 1), \\ & \quad (S_{3,5,1}, (0, 2, 0), 1), \\ & \quad (S_{4,1,1}, (1, -1, 0), 1), \\ (S_{4,2,1}, (1, 0, -1), 1), & (S_{4,2,2}, (1, 0, 0), 0), (S_{4,2,3}, (1, 0, 1), 1), \\ & \quad (S_{4,3,1}, (1, 1, 0), 1), \\ & \quad (S_{5,1,1}, (2, 0, 0), 1)\}. \end{aligned}$$

This example is particularly simple because we can choose all sample points with rational coordinates. This will not be the case in general: the coordinates of the sample points will be roots of univariate polynomials above sample points of cells of lower dimension, and the real roots techniques of Chapter 10 will have to be generalized to deal with the cylindrical situation.

### 11.1.2 Details of the Lifting Phase

In order to make precise the lifting phase of the Cylindrical Decomposition Algorithm, it is necessary to compute sample points on cells. In the archimedean case, this can be done using isolating intervals. For a general real closed field, Thom encodings will be used.

Since the degree bounds are already doubly exponential, and since we are going to give a much better algorithm in Chapter 14, we do not dwell on precise complexity analysis of the lifting phase.

The notion of a triangular system of equations is natural in the context of the lifting phase of the Cylindrical Decomposition Algorithm.

**Definition 11.3.** Let  $K$  be a field contained in an algebraically closed field  $C$ . A **triangular system of polynomials** with variables  $X_1, \dots, X_k$  is a list  $\mathcal{T} = \mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k$ , where

$$\begin{aligned} \mathcal{T}_1 &\in K[X_1], \\ \mathcal{T}_2 &\in K[X_1, X_2], \\ &\vdots \\ \mathcal{T}_k &\in K[X_1, \dots, X_k], \end{aligned}$$

such that the polynomial system  $\mathcal{T}$  is zero-dimensional, i.e.  $\text{Zer}(\mathcal{T}, C^k)$  is finite. □

### 11.1.2.1 The Archimedean Case

In this case, we are going to use isolating intervals to characterize the sample points of the cells. We need a notion of parallelepiped isolating a point.

A **parallelepiped isolating**  $z \in \mathbb{R}^k$  is a list  $(\mathcal{T}_1, I_1), (\mathcal{T}_2, I_2), \dots, (\mathcal{T}_k, I_k)$  where  $\mathcal{T}_1 \in \mathbb{R}[X_1], \dots, \mathcal{T}_k \in \mathbb{R}[X_1, \dots, X_k], \mathcal{T} = \mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k$  is a triangular system,  $I_1$  is an open interval with rational end points or a rational containing the root  $z_1$  of  $\mathcal{T}_1$  and no other root of  $\mathcal{T}_1$  in  $\mathbb{R}$ ,  $I_2$  is an open interval with rational end points or a rational containing the root  $z_2$  of  $\mathcal{T}_2(z_1, X_2)$  and no other root of  $\mathcal{T}_2(z_1, X_2)$  in  $\mathbb{R}$ , ...,  $I_k$  is an open interval with rational end points or a rational containing the root  $z_k$  of  $\mathcal{T}_k(z_1, \dots, z_{k-1}, X_k)$  and no other root of  $\mathcal{T}_k(z_1, \dots, z_{k-1}, X_k)$  in  $\mathbb{R}$ .

Given a parallelepiped isolating  $z \in \mathbb{R}^k$  it is not difficult, using elementary properties of intervals, to give bounds on the value of  $P(z)$  where  $P$  is a polynomial of  $\mathbb{R}[X_1, \dots, X_k]$ . So if  $Q \in \mathbb{R}[X_1, \dots, X_k, X_{k+1}]$ , it is not difficult to find a natural number  $N$  such that all the roots of  $Q(z, X_k)$  belong to  $(-2^N, 2^N)$  using Lemma 10.6.

As in the case of a univariate polynomial, the root isolation method can be used for characterizing real roots in the cylindrical situation. Note that testing equality to zero and deciding signs are necessary to evaluate the degrees of the polynomials and the numbers of sign variations in their coefficients. These testing equality to zero and deciding signs will be done through a recursive call to Algorithm 11.4 (Multivariate Sign at a Sample Point).

Let us first consider an example in order to illustrate this situation in the simple situation where  $k = 2$ .

*Example 11.4.* We want to isolate the real roots of the polynomials

$$\mathcal{T}_1 = 9X_1^2 - 1, \mathcal{T}_2 = (5X_2 - 1)^2 + 3X_1 - 1.$$

We first isolate the roots of  $\mathcal{T}_1$  and get  $z_1$  isolated by  $(9X_1^2 - 1, [0, 1])$  and  $z'_1$  isolated by  $(9X_1^2 - 1, [-1, 0])$ .

We now want to isolate the roots of  $\mathcal{T}_2(z_1, X_2) = (5X_2 - 1)^2 + 3z_1 - 1$ , using Algorithm 10.4 (Real Root Isolation) in a recursive way. We first need to know the separable part of  $\mathcal{T}_2(z_1, X_2)$ . In order to compute it, it is necessary to decide whether  $3z_1 - 1$  is 0. For this we call Algorithm 10.6 which computes the gcd of  $3X_1 - 1$  and  $9X_1^2 - 1$ , which is  $3X_1 - 1$ , and checks whether  $3X_1 - 1$  vanishes at  $z_1$ . This is the case since the sign of  $3X_1 - 1$  changes between 0 and 1. So the separable part of  $\mathcal{T}_2(z_1, X_2)$  is  $5X_2 - 1$  and  $\mathcal{T}_2(z_1, X_2)$  has a single double root above  $z_1$ .

We now isolate the roots of  $\mathcal{T}_2(z'_1, X_2) = (5X_2 - 1)^2 + 3z'_1 - 1$ . We follow again the method of Algorithm 10.4 (Real Root Isolation). We first need to know the separable part of  $\mathcal{T}_2(z'_1, X_2)$ . In order it, it is necessary to decide whether  $3z'_1 - 1$  is 0. For this purpose we call Algorithm 10.6 which computes the gcd of  $3X_1 - 1$  and  $9X_1^2 - 1$ , which is  $3X_1 - 1$ , and checks whether  $3X_1 - 1$  vanishes at  $z'_1$ . This is not the case, since the sign of  $3X_1 - 1$  does not changes between -1 and 0. In fact,  $3z'_1 - 1$  is negative. So  $\mathcal{T}_2(z'_1, X_2)$  is separable.

Continuing the isolation process, we finally find that  $P_2(z'_1, X_2)$  has two distinct real roots, one positive and one negative.  $\square$

Note that in this example it was not necessary to refine the intervals defining  $z_1$  and  $z'_1$  to decide whether  $\mathcal{T}_2(z_1, X_2)$  and  $\mathcal{T}_2(z'_1, X_2)$  were separable. However, in the general situation considered now, such refinements may be necessary, and are produced by the recursive calls.

### Algorithm 11.3. [Real Recursive Root Isolation]

- **Structure:** the field of real numbers  $\mathbb{R}$ .
- **Input:** a parallelepiped isolating  $z \in \mathbb{R}^k$ , a polynomial  $P \in \mathbb{R}[X_1, \dots, X_{k+1}]$ , and a natural number  $N$  such that all the roots of  $P(z, X_k)$  belong to  $(-2^N, 2^N)$ .
- **Output:** a parallelepipeds isolating  $(z, y)$  for every  $y$  root of  $P(z, X_{k+1})$ .
- **Procedure:** Perform the computations in Algorithm 10.4 (Real Root Isolation), testing equality to zero and deciding signs necessary for the computation of the degrees and of the sign variations being done by recursive calls to Algorithm 11.4 (Real Recursive Sign at a Point) at level  $k$ .

So we need to find the sign of a polynomial at a point. This algorithm calls itself recursively.

### Algorithm 11.4. [Real Recursive Sign at a Point]

- **Structure:** the field of real numbers  $\mathbb{R}$ .
- **Input:** a parallelepiped isolating  $z \in \mathbb{R}^k$ , a polynomial  $Q(X_1, \dots, X_k)$ , and a natural number  $N$  such that all the roots of  $Q(z_1, \dots, z_{k-1}, X_k)$  belong to  $(-2^N, 2^N)$ .
- **Output:** a parallelepiped isolating  $z \in \mathbb{R}^k$ , and the sign of  $Q(z)$ .
- **Procedure:**
  - If  $k = 1$ , perform Algorithm 10.6 (Sign at a Real Root).

- If  $k > 1$ , perform the computations of the Algorithm 10.6 (Sign at a Real Root), testing equality to zero and deciding signs necessary for the computation of the degrees and of the sign variations being done by recursive calls to Algorithm 11.4 (Multivariate Sign at a Point) with level  $k - 1$ .

We can compare the roots of two polynomials. Again Algorithm 11.4 is called to evaluate sign variations.

*Algorithm 11.5. [Recursive Comparison of Real Roots]*

- **Structure:** the field of real numbers  $\mathbb{R}$ .
- **Input:** a parallelepiped isolating  $z \in \mathbb{R}^{k-1}$ , a polynomial  $P(X_1, \dots, X_k)$ , a polynomial  $Q(X_1, \dots, X_k)$ , and a natural number  $N$  such that all the roots of  $P(z, X_k)$ ,  $Q(z, X_k)$  belong to  $(-2^N, 2^N)$ .
- **Output:** a parallelepiped isolating  $(z, y)$  for every root  $y$  of  $P(z, X_k)$  or  $Q(z, X_k)$ , and the signs of  $Q(z, y)$  (resp.  $P(z, y)$ ).
- **Procedure:** Perform the computations of the Algorithm 10.7 (Comparison of Real Roots), testing equality to zero and deciding signs necessary for the computation of the degrees and of the sign variations being done by recursive calls to Algorithm 11.4 (Real Recursive Sign at a Point) with level  $k - 1$ .

Finally, we can find sample points for a family of polynomials above a point.

*Algorithm 11.6. [Real Recursive Sample Points]*

- **Structure:** the field of real numbers  $\mathbb{R}$ .
- **Input:** a parallelepiped isolating  $z \in \mathbb{R}^{k-1}$ , a finite set of polynomials  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ , and a natural number  $N$  such that all the roots of  $P(z, X_k)$  belong to  $(-2^N, 2^N)$  for  $P \in \mathcal{P}$ .
- **Output:** a level  $k$ , parallelepipeds isolating the roots of the non-zero polynomials in  $\mathcal{P}(z, X_k)$ , an element between two consecutive roots of elements of  $\mathcal{P}(z, X_k)$ , an element of  $\mathbb{R}$  smaller than all these roots and an element of  $\mathbb{R}$  greater than all these roots. The sign of all  $Q(z, y)$ ,  $Q \in \mathcal{P}$  is also output for every root of an element of  $\mathcal{P}(z, X_k)$ .
- **Procedure:**
  - For every pair  $P, Q$  of elements of  $\mathcal{P}$ , perform Algorithm 11.5.
  - Compute a rational point in between two consecutive roots using the isolating intervals.
  - Compute a rational point smaller than all these roots and rational point greater than all the roots of polynomials in  $\mathcal{P}(z, X_k)$  using Proposition 10.1.

The preceding algorithms make it possible to describe more precisely the Lifting Phase of the Cylindrical Decomposition Algorithm 11.2 in the real case.

**Algorithm 11.7. [Real Lifting Phase]**

- **Structure:** the field of real numbers  $\mathbb{R}$ .
- **Input:**  $C_i(\mathcal{P})$ , for  $i = k - 1, \dots, 1$ .
- **Output:** a cylindrical set of sample points of a cylindrical decomposition  $\mathcal{S}_1, \dots, \mathcal{S}_k$  of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$  and, for each sample point, the sign of the polynomials in  $\mathcal{P}$  at this point.
- **Procedure:**
  - Run Algorithm 10.8 (Sample Points on a Line) with input  $C_1(\mathcal{P})$  to obtain the sample points of the cells in  $\mathcal{S}_1$ .
  - For every  $i = 2, \dots, k$ , compute the sample points of the cells of  $\mathcal{S}_i$  from the sample points of the cells in  $\mathcal{S}_{i-1}$  as follows: Compute for every parallelepiped  $I$  specifying  $x$  a list denoted by  $L$  of non-zero polynomials  $P_i(x, X_i)$  with  $P_i \in C_i(\mathcal{P})$  using Algorithm 11.4 (Real Recursive Sign at a Point). Run Algorithm 11.6 (Real Recursive Sample Points) with input  $L$  and  $x$ .

**11.1.2.2 The case of a real closed field**

In this case, we are going to use Thom encodings to characterize the sample points of the cells.

**Definition 11.5. A triangular Thom encoding specifying**

$$z = (z_1, \dots, z_k) \in \mathbb{R}^k$$

is a pair  $\mathcal{T}, \sigma$  where  $\mathcal{T}$  is a triangular system of polynomials and  $\sigma = \sigma_1, \dots, \sigma_k$  is a list of Thom encodings such that

- $\sigma_1$  is the Thom encoding of a root  $z_1$  of  $\mathcal{T}_1$ ,
- $\sigma_2$  is the Thom encoding of a root  $z_2$  of  $\mathcal{T}_2(z_1, X_2)$ ,
- $\vdots$
- $\sigma_k$  is the Thom encoding of a root  $z_k$  of  $\mathcal{T}_k(z_1, \dots, z_{k-1}, X_k)$ .

In other words, denoting by  $\text{Der}(\mathcal{T})$  the set of derivatives of  $\mathcal{T}_j$  with respect to  $X_j$ ,  $j = 1, \dots, k$ ,  $\sigma$  is a sign condition on  $\text{Der}(\mathcal{T})$ .

We denote by  $z(\mathcal{T}, \sigma)$  the  $k$ -tuple specified by the Thom encoding  $\mathcal{T}, \sigma$ .  $\square$

The lifting phase of the Cylindrical Decomposition Algorithm in the case of a real closed field is based on the following recursive algorithms generalizing the corresponding univariate algorithms in Chapter 10.

**Algorithm 11.8. [Recursive Sign Determination]**

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $K$ .
- **Input:** a triangular system  $\mathcal{T}$ , and a list  $\mathcal{Q}$  of elements of  $D[X_1, \dots, X_k]$ . Denote by  $Z = \text{Zer}(\mathcal{T}, \mathbb{R}^k)$ .



- **Output:** the set  $\text{SIGN}(\mathcal{Q}, \mathcal{T})$  of sign conditions realized by  $\mathcal{Q}$  on  $Z$ .
- **Procedure:**
  - If  $k = 1$ , perform Algorithm 10.13 (Univariate Sign Determination).
  - If  $k > 1$ , perform the computations of the Algorithm 10.11 (Sign Determination), deciding signs necessary for the determination of the necessary Tarski-queries by recursive calls to Algorithm 11.8 (Recursive Sign Determination) with level  $k - 1$ .

**Exercise 11.1.** Prove that the complexity of Algorithm 11.8 (Recursive Sign Determination) is  $s d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{Q}$ ,  $d$  is a bound on the degrees on the elements of  $\mathcal{T}$  and  $\mathcal{Q}$ .

*Algorithm 11.9.* **[Recursive Thom Encoding]**

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a triangular system  $\mathcal{T}$ .
- **Output:** the list  $\text{Thom}(\mathcal{T})$  of Thom encodings of the roots of  $\mathcal{T}$ .
- **Procedure:** Apply Algorithm 11.8 (Recursive Sign Determination) to  $\mathcal{T}$  and  $\text{Der}(\mathcal{T})$ .

**Exercise 11.2.** Prove that the complexity of Algorithm 11.9 (Recursive Thom Encoding) is  $d^{O(k)}$ , where  $d$  is a bound on the degrees on the elements of  $\mathcal{T}$ .

Let  $\mathcal{T}, \sigma$  be a triangular Thom encoding specifying a point  $z = (z_1, \dots, z_{k-1})$  of  $R^{k-1}$ .

**Definition 11.6.** A Thom encoding  $P, \tau$  above  $\mathcal{T}, \sigma$  is

- a polynomial  $P \in R[X_1, \dots, X_k]$ ,
- a sign condition  $\tau$  on  $\text{Der}_{X_k}(P)$  such that  $\sigma, \tau$  is the triangular Thom encoding of a root  $(z, a)$  of  $\mathcal{T}, P$ . □

*Algorithm 11.10.* **[Recursive Comparison of Roots]**

- **Structure:** an ordered integral domain  $D$ , contained in a real closed field  $R$ .
- **Input:** a Thom encoding  $\mathcal{T}, \sigma$  specifying  $z \in R^{k-1}$ , and two non-zero polynomials  $P$  and  $Q$  in  $D[X_1, \dots, X_k]$ .
- **Output:** the ordered list of the Thom encodings of the roots of  $P$  and  $Q$  over  $\sigma$ .
- **Procedure:** Apply Algorithm 11.8 (Recursive Sign Determination)

$$\mathcal{T}, P, \text{Der}(\mathcal{T}) \cup \text{Der}(P) \cup \text{Der}(Q),$$

then to

$$\mathcal{T}, Q, \text{Der}(\mathcal{T}) \cup \text{Der}(Q) \cup \text{Der}(P).$$

Compare the roots using Proposition 2.28.

**Exercise 11.3.** Prove that the complexity of Algorithm 11.10 (Recursive Comparison of Roots) is  $d^{O(k)}$ , where  $d$  is a bound on the degrees on the elements of  $\mathcal{T}$ ,  $P$  and  $Q$ .

We can also construct points between two consecutive roots.

*Algorithm 11.11. [Recursive Intermediate Points]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a Thom encoding  $\mathcal{T}$ ,  $\sigma$  specifying  $z \in R^{k-1}$ , and two non-zero polynomials  $P$  and  $Q$  in  $D[X_1, \dots, X_k]$ .
- **Output:** Thom encodings specifying values  $y$  intersecting intervals between two consecutive roots of  $P(z, X_k)$  and  $Q(z, X_k)$ .
- **Procedure:** Compute the Thom encodings of the roots of the polynomial  $\partial(PQ)/\partial X_k(z, X_k)$  above  $\mathcal{T}$ ,  $\sigma$  using Algorithm 11.9 (Recursive Thom Encoding) and compare them to the roots of  $P$  and  $Q$  above  $\sigma$  using Algorithm 11.10 (Recursive Comparison of Roots). Keep one intermediate point between two consecutive roots of  $PQ$ .

**Exercise 11.4.** Prove that the complexity of Algorithm 11.11 (Recursive Intermediate Points) is  $d^{O(k)}$ , where  $d$  is a bound on the degrees on the elements of  $\mathcal{T}$ ,  $P$  and  $Q$ .

Finally we can compute sample points on a line.

*Algorithm 11.12. [Recursive Sample Points]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a Thom encoding  $\mathcal{T}$ ,  $\sigma$  specifying  $z \in R^{k-1}$ , and a family of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** an ordered list  $L$  of Thom encodings specifying the roots in  $R$  of the non-zero polynomials  $P(z, X_k)$ ,  $P \in \mathcal{P}$ , an element between two such consecutive roots, an element of  $R$  smaller than all these roots, and an element of  $R$  greater than all these roots. Moreover  $(\tau_1)$  appears before  $(\tau_2)$  in  $L$  if and only if  $x_k(\tau_1) \leq x_k(\tau_2)$ . The sign of  $Q(z, x_k(\tau))$  is also output for every  $Q \in \mathcal{P}$ ,  $\tau \in L$ .
- **Procedure:** Characterize the roots of the polynomials in  $R$  using Algorithm 11.9 (Recursive Thom Encoding). Compare these roots using Algorithm 11.10 (Recursive Comparison of Roots) for every pair of polynomials in  $\mathcal{P}$ . Characterize a point in each interval between the roots by Algorithm 11.11 (Recursive Intermediate Points). Use Proposition 10.5 to find an element of  $R$  smaller and bigger than any root of any polynomial in  $\mathcal{P}$  above  $z$ .

**Exercise 11.5.** Prove that the complexity of Algorithm 11.12 (Recursive Sample Points) is  $s d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{Q}$  and  $d$  is a bound on the degrees on the elements of  $\mathcal{T}$  and  $\mathcal{Q}$ .

We are now ready to describe the lifting phase of the Cylindrical Decomposition Algorithm in the general case.

*Algorithm 11.13. [Lifting Phase]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:**  $C_i(\mathcal{P})$ , for  $i = k, \dots, 1$ .
- **Output:** a set of sample points of a cylindrical decomposition  $\mathcal{S}_1, \dots, \mathcal{S}_k$ , of  $R^k$  adapted to  $\mathcal{P}$  and for each sample point the sign of the polynomials in  $\mathcal{P}$  at this point.
- **Procedure:**
  - Run Algorithm 10.19 (Cylindrical Univariate Sample Points) with input  $C_1(\mathcal{P})$  to obtain the sample points of the cells in  $\mathcal{S}_1$ , (described by Thom encodings).
  - For every  $i = 1, \dots, k - 1$ , compute the sample points of the cells of  $\mathcal{S}_i$  from the sample points of the cells in  $\mathcal{S}_{i-1}$ , (described by triangular Thom encodings) as follows: Compute for every  $\sigma$  specifying a sample point  $x$  the list  $L$  of non-zero polynomials  $P_i(x, X_i)$  with  $P_i \in C_i(\mathcal{P})$ , using Algorithm 11.8 (Recursive Sign Determination). Run Algorithm 11.12 (Recursive Sample Points) with input  $L$  and  $\sigma$ .

**Exercise 11.6.** Prove that the complexity of Algorithm 11.13 (Lifting Phase) is  $(s d)^{O(1)^k}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{P}$ .

## 11.2 Decision Problem

Now we explain how to decide the truth or falsity of a sentence using the Cylindrical Decomposition Algorithm applied to the family of polynomials used to build the sentence.

Let  $\mathcal{P}$  be a finite subset of  $R[X_1, \dots, X_k]$ . A  **$\mathcal{P}$ -atom** is one of  $P = 0$ ,  $P \neq 0$ ,  $P > 0$ ,  $P < 0$ , where  $P$  is a polynomial in  $\mathcal{P}$ . A  **$\mathcal{P}$ -formula** is a formula (Definition page 58) written with  $\mathcal{P}$ -atoms. A  **$\mathcal{P}$ -sentence** is a sentence (Definition page 59) written with  $\mathcal{P}$ -atoms.

**Notation 11.7. [Cylindrical realizable sign conditions]**

For  $z \in R^k$ , we denote by  $\text{sign}(\mathcal{P})(z)$  the sign condition on  $\mathcal{P}$  mapping  $P \in \mathcal{P}$  to  $\text{sign}(P)(z) \in \{0, 1, -1\}$ .

We are going to define inductively the tree of cylindrical realizable sign conditions,  $\text{CSIGN}(\mathcal{P})$ , of  $\mathcal{P}$ . The importance of this notion is that the truth or falsity of any  $\mathcal{P}$ -sentence can be decided from  $\text{CSIGN}(\mathcal{P})$ .

We denote by  $\pi_i$  the projection from  $\mathbb{R}^{i+1}$  to  $\mathbb{R}^i$  forgetting the last coordinate. By convention,  $\mathbb{R}^0 = \{0\}$ .

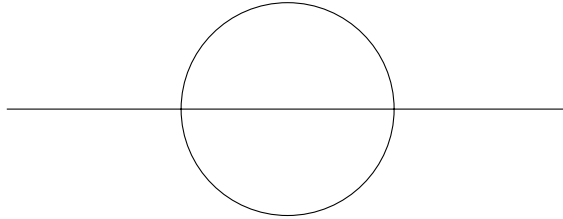
- For  $z \in \mathbb{R}^k$ , let  $\text{CSIGN}_k(\mathcal{P})(z) = \text{sign}(\mathcal{P})(z)$ .
- For  $i, 0 \leq i < k$ , and all  $y \in \mathbb{R}^i$ , we inductively define

$$\text{CSIGN}_i(\mathcal{P})(y) = \{\text{CSIGN}_{i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{i+1}, \pi_i(z) = y\}.$$

Finally, we define the **tree of cylindrical realizable sign conditions of  $\mathcal{P}$** , denoted  $\text{CSIGN}(\mathcal{P})$ , by

$$\text{CSIGN}(\mathcal{P}) = \text{CSIGN}_0(\mathcal{P})(0). \quad \square$$

*Example 11.8.* Consider two bivariate polynomials  $P_1 = X_2, P_2 = X_1^2 + X_2^2 - 1$  and  $\mathcal{P} = \{P_1, P_2\}$ .



**Fig. 11.1.**  $\text{Zer}(P_1, \mathbb{R}^2)$  and  $\text{Zer}(P_2, \mathbb{R}^2)$

We order the set  $\mathcal{P}$  with the order  $P_1 < P_2$ .

For  $y \in \mathbb{R}^2$ ,  $\text{sign}(\mathcal{P})(y)$  is the mapping from  $\mathcal{P}$  to  $\{0, 1, -1\}$  sending  $(P_1, P_2)$  to  $(\text{sign}(P_1(y)), \text{sign}(P_2(y)))$ . Abusing notation, we denote the mapping  $\text{sign}(\mathcal{P})(y)$  by  $(\text{sign}(P_1(y)), \text{sign}(P_2(y)))$ .

For example if  $y = (0, 0)$ ,  $\text{sign}(\mathcal{P})(0, 0) = (0, -1)$  since  $\text{sign}(P_1(0, 0)) = 0$  and  $\text{sign}(P_2(0, 0)) = -1$ .

Fixing  $x \in \mathbb{R}$ ,  $\text{CSIGN}_1(\mathcal{P})(x)$  is the set of all possible  $\text{sign}(\mathcal{P})(z)$  for  $z \in \mathbb{R}^2$  such that  $\pi_1(z) = x$ . For example if  $x = 0$ , there are seven possibilities for  $\text{sign}(\mathcal{P})(z)$  as  $z$  varies in  $\{0\} \times \mathbb{R}$ :

$$(-1, 1), (-1, 0), (-1, -1), (0, -1), (1, -1), (1, 0), (1, 1).$$

So  $\text{CSIGN}_1(\mathcal{P})(0)$  is

$$\{(-1, 1), (-1, 0), (-1, -1), (0, -1), (1, -1), (1, 0), (1, 1)\}.$$

Similarly, if  $x = 1$ , there are three possibilities for  $\text{sign}(\mathcal{P})(z)$  as  $z$  varies in  $\{1\} \times \mathbb{R}$ :

$$(-1, 1), (0, 0), (1, 1).$$

So  $\text{CSIGN}_1(\mathcal{P})(1)$  is

$$\{(-1, 1), (0, 0), (1, 1)\}.$$

If  $x = 2$ , there are three possibilities for  $\text{sign}(\mathcal{P})(z)$  as  $z$  varies in  $\{2\} \times \mathbb{R}$ :

$$(-1, 1), (0, 1), (1, 1).$$

So  $\text{CSIGN}_1(\mathcal{P})(2)$  is

$$\{(-1, 1), (0, 1), (1, 1)\}.$$

Finally  $\text{CSIGN}(\mathcal{P})$  is the set of all possible  $\text{CSIGN}_1(\mathcal{P})(x)$  for  $x \in \mathbb{R}$ . It is easy to check that the three cases we have considered ( $x = 0, x = 1, x = 2$ ) already give all possible  $\text{CSIGN}_1(\mathcal{P})(x)$  for  $x \in \mathbb{R}$ . So  $\text{CSIGN}(\mathcal{P})$  is the set with three elements

$$\begin{aligned} & \{ \{(-1, 1), (-1, 0), (-1, -1), (0, -1), (1, -1), (1, 0), (1, 1)\}, \\ & \quad \{(-1, 1), (0, 0), (1, 1)\}, \\ & \quad \{(-1, 1), (0, 1), (1, 1)\} \}. \end{aligned} \quad \square$$

We now explain how  $\text{CSIGN}(\mathcal{P})$  can be determined from a cylindrical set of sample points of a cylindrical decomposition adapted to  $\mathcal{P}$  and the signs of  $P \in \mathcal{P}$  at these points.

If  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_k, \mathcal{A}_i \subset \mathbb{R}^k, \pi_i(\mathcal{A}_{i+1}) = \mathcal{A}_i$ , where  $\pi_i$  is the projection from  $\mathbb{R}^{i+1}$  to  $\mathbb{R}^i$  forgetting the last coordinate, we define inductively the **tree of cylindrical realizable sign conditions**  $\text{CSIGN}(\mathcal{P}, \mathcal{A})$  of  $\mathcal{P}$  on  $\mathcal{A}$ .

– For  $z \in \mathcal{A}_k$ , let

$$\text{CSIGN}_k(\mathcal{P}, \mathcal{A})(z) = \text{sign}(\mathcal{P})(z).$$

– For all  $i, 0 \leq i < k$ , and all  $y \in \mathcal{A}_i$ , we inductively define

$$\text{CSIGN}_i(\mathcal{P}, \mathcal{A})(y) = \{ \text{CSIGN}_{i+1}(\mathcal{P}, \mathcal{A})(z) \mid z \in \mathcal{A}_{i+1}, \pi_i(z) = y \}.$$

Finally,

$$\text{CSIGN}(\mathcal{P}, \mathcal{A}) = \text{CSIGN}_0(\mathcal{P}, \mathcal{A})(0).$$

Note that  $\text{CSIGN}(\mathcal{P}) = \text{CSIGN}(\mathcal{P}, \mathbb{R}^k)$ . Note also that  $\text{CSIGN}(\mathcal{P}, \mathcal{A})$  is a subtree of  $\text{CSIGN}(\mathcal{P})$ .

We are going to prove the following result.

**Proposition 11.9.** *Let  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_k$  be a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$  and let  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_k$  be a cylindrical set of sample points for  $\mathcal{S}$ . Then*

$$\text{CSIGN}(\mathcal{P}, \mathcal{A}) = \text{CSIGN}(\mathcal{P}).$$

We first start by explaining how this works on an example.

*Example 11.10.* Let  $P = X_1^2 + X_2^2 + X_3^2 - 1$  and  $\mathcal{P} = \{P\}$ . Since there is only one polynomial in  $\mathcal{P}$ , we identify  $\{0, 1, -1\}^{\mathcal{P}}$  with  $\{0, 1, -1\}$ .

We use Example 11.2, where the cells and sample points of the cylindrical decomposition of  $\{P = X_1^2 + X_2^2 + X_3^2 - 1\}$  were described.

The sign condition  $\text{sign}(\mathcal{P})(u)$  is fixed on each cell of  $\mathbb{R}^3$  by the sign of  $P$  at the sample point of the cell and thus

$$\text{sign}(\mathcal{P})(z) = \begin{cases} -1 & \text{if } z \in S_{3,3,3} \\ 0 & \text{if } z \in S_{2,2,1} \cup S_{2,2,2} \cup S_{3,2,2} \\ & \cup S_{3,3,2} \cup S_{3,3,4} \cup S_{3,4,2} \cup S_{4,2,2} \\ 1 & \text{otherwise.} \end{cases}$$

The set  $\text{CSIGN}_2(\mathcal{P})(y)$  is fixed on each cell of  $\mathbb{R}^2$  by its value at the sample point of the cell and thus

$$\text{CSIGN}_2(\mathcal{P})(y) = \begin{cases} \{0, 1, -1\} & \text{if } y \in S_{3,3} \\ \{0, 1\} & \text{if } y \in S_{2,2} \cup S_{3,2} \cup S_{3,4} \cup S_{4,2} \\ \{1\} & \text{otherwise.} \end{cases}$$

The set  $\text{CSIGN}_1(\mathcal{P})(x)$  is fixed on each cell of  $\mathbb{R}$  by its value at the sample point of the cell and thus

$$\text{CSIGN}_1(\mathcal{P})(x) = \begin{cases} \{\{1\}, \{0, 1\}, \{0, 1, -1\}\} & \text{if } x \in S_3 \\ \{\{1\}, \{0, 1\}\} & \text{if } x \in S_2 \cup S_4 \\ \{\{1\}\} & \text{if } x \in S_1 \cup S_5. \end{cases}$$

Finally the set  $\text{CSIGN}(\mathcal{P})$  has three elements and

$$\text{CSIGN}(\mathcal{P}) = \{\{\{1\}, \{0, 1\}, \{0, 1, -1\}\}, \{\{1\}, \{0, 1\}\}, \{\{1\}\}\}.$$

This means that there are three possible cases:

- there are values of  $x_1 \in \mathbb{R}$  for which
  - for some value of  $x_2 \in \mathbb{R}$ , the only sign taken by  $P(x_1, x_2, x_3)$  when  $x_3$  varies in  $\mathbb{R}$  is 1,
  - for some value of  $x_2 \in \mathbb{R}$ , the only signs taken by  $P(x_1, x_2, x_3)$  when  $x_3$  varies in  $\mathbb{R}$  are 0 or 1,
  - for some value of  $x_2 \in \mathbb{R}$ , the signs taken by  $P(x_1, x_2, x_3)$  when  $x_3$  varies in  $\mathbb{R}$  are 0, 1, or  $-1$ ,
  - and these are the only possibilities,
- there are values of  $x_1$  for which
  - for some value of  $x_2 \in \mathbb{R}$ , the only sign taken by  $P(x_1, x_2, x_3)$  when  $x_3$  varies in  $\mathbb{R}$  is 1,
  - for some value of  $x_2 \in \mathbb{R}$ , the only signs taken by  $P(x_1, x_2, x_3)$  when  $x_3$  varies in  $\mathbb{R}$  are 0 or 1,
  - and these are the only possibilities,
- there are values of  $x_1$  for which
  - the only sign taken by  $P(x_1, x_2, x_3)$  when  $(x_2, x_3)$  varies in  $\mathbb{R}^2$  is 1,

– and together these three cases exhaust all possible values of  $x_1 \in \mathbb{R}$ .  $\square$

**Proposition 11.11.** *Let  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_k$  be a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ . For every  $1 \leq i \leq k$  and every  $S \in \mathcal{S}_i$ ,  $\text{CSIGN}_i(y)$  is constant as  $y$  varies in  $S$ .*

**Proof:** The proof is by induction on  $k - i$ .

If  $i = k$ , the claim is true since the sign of every  $P \in \mathcal{P}$  is fixed on  $S \in \mathcal{S}_k$ .

Suppose that the claim is true for  $i + 1$  and consider  $S \in \mathcal{S}_i$ . Let  $T_1, \dots, T_\ell$  be the cells of  $\mathcal{S}_{i+1}$  such that  $\pi_i(T_j) = S$ . By induction hypothesis,  $\text{CSIGN}_{i+1}(\mathcal{P})(z)$  is constant as  $z$  varies in  $T_j$ . Since  $\mathcal{S}$  is a cylindrical decomposition,  $\bigcup_{j=1}^\ell T_j = S \times \mathbb{R}$ . Thus

$$\text{CSIGN}_i(\mathcal{P})(y) = \{\text{CSIGN}_{i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{i+1}, \pi_i(z) = y\}$$

is constant as  $y$  varies in  $S$ .  $\square$

**Proof of Proposition 11.9:** Let  $\mathcal{A}_0 = \{0\}$ . We are going to prove that for every  $y \in \mathcal{A}_i$ ,

$$\text{CSIGN}_i(\mathcal{P})(y) = \text{CSIGN}_i(\mathcal{P}, \mathcal{A})(y).$$

The proof is by induction on  $k - i$ .

If  $i = k$ , the claim is true since  $\mathcal{A}_k$  meets every cell of  $\mathcal{S}_k$ .

Suppose that the claim is true for  $i + 1$  and consider  $y \in \mathcal{A}_i$ . Let  $S \in \mathcal{S}_i$  be the cell containing  $y$ , and let  $T_1, \dots, T_\ell$  be the cells of  $\mathcal{S}_{i+1}$  such that  $\pi_i(T_j) = S$ . Denote by  $z_j$  the unique point of  $T_j \cap \mathcal{A}_{i+1}$  such that  $\pi_i(z_j) = y$ . By induction hypothesis,

$$\text{CSIGN}_{i+1}(\mathcal{P})(z_j) = \text{CSIGN}_{i+1}(\mathcal{P}, \mathcal{A})(z_j).$$

Since  $\text{CSIGN}_{i+1}(\mathcal{P})(z)$  is constant as  $z$  varies in  $T_j$ ,

$$\begin{aligned} \text{CSIGN}_i(\mathcal{P})(y) &= \{\text{CSIGN}_{i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{i+1}, \pi_i(z) = y\} \\ &= \{\text{CSIGN}_{i+1}(\mathcal{P}, \mathcal{A})(z) \mid z \in \mathcal{A}_{i+1}, \pi_i(z) = y\} \\ &= \text{CSIGN}_i(\mathcal{P}, \mathcal{A})(y) \end{aligned}$$

$\square$

The Cylindrical Decision Algorithm is based on the following result. We are going to need a notation.

**Notation 11.12.** If  $\mathcal{P} \subset \mathbb{K}[X_1, \dots, X_k]$  is finite,  $X = (X_1, \dots, X_k)$ ,  $F(X)$  is a  $\mathcal{P}$ -quantifier free formula, and  $\sigma \in \mathcal{P}^{\{0,1,-1\}}$  is a sign condition on  $\mathcal{P}$ , we define  $F^*(\sigma) \in \{\text{True}, \text{False}\}$  as follows:

– If  $F$  is the atom  $P = 0$ ,  $P \in \mathcal{P}$ ,  $F^*(\sigma) = \text{True}$  if  $\sigma(P) = 0$ ,  $F^*(\sigma) = \text{False}$  otherwise.

- If  $F$  is the atom  $P > 0$ ,  $P \in \mathcal{P}$ ,  $F^*(\sigma) = \text{True}$  if  $\sigma(P) = 1$ ,  $F^*(\sigma) = \text{False}$  otherwise.
- If  $F$  is the atom  $P < 0$ ,  $P \in \mathcal{P}$ ,  $F^*(\sigma) = \text{True}$  if  $\sigma(P) = -1$ ,  $F^*(\sigma) = \text{False}$  otherwise.
- If  $F = F_1 \wedge F_2$ ,  $F^*(\sigma) = F_1^*(\sigma) \wedge F_2^*(\sigma)$ .
- If  $F = F_1 \vee F_2$ ,  $F^*(\sigma) = F_1^*(\sigma) \vee F_2^*(\sigma)$ .
- If  $F = \neg(G)$ ,  $F^*(\sigma) = \neg(G^*(\sigma))$ . □

*Example 11.13.* If  $F = X_1^2 + X_2^2 + X_3^2 - 1 > 0$ , then

$$F^*(\sigma) = \begin{cases} \text{True} & \text{if } \sigma = 1 \\ \text{False} & \text{if } \sigma = 0, -1 \end{cases} \quad \square$$

**Proposition 11.14.** *The  $\mathcal{P}$ -sentence*

$$(\text{Qu}_1 X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(X_1, \dots, X_k),$$

where  $F(X_1, \dots, X_k)$  is quantifier free,  $\text{Qu}_i \in \{\exists, \forall\}$ , is true if and only if

$$(\text{Qu}_1 \sigma_1 \in \text{CSIGN}(\mathcal{P})) (\text{Qu}_2 \sigma_2 \in \sigma_1) \dots (\text{Qu}_k \sigma_k \in \sigma_{k-1}) F^*(\sigma_k)$$

is true.

*Example 11.15.* We illustrate the statement of the proposition by an example. Consider again  $\mathcal{P} = \{X_1^2 + X_2^2 + X_3^2 - 1\}$ , and recall that

$$\text{CSIGN}(\mathcal{P}) = \{\{\{1\}, \{0, 1\}, \{0, 1, -1\}\}, \{\{1\}, \{0, 1\}\}, \{\{1\}\}\}$$

by Example 11.2.

The sentence  $(\forall X_1)(\forall X_2)(\forall X_3) F$ , with  $F = X_1^2 + X_2^2 + X_3^2 - 1 > 0$  is false since taking  $(x_1, x_2, x_3) = (0, 0, 0)$  we get  $x_1^2 + x_2^2 + x_3^2 - 1 < 0$ . It is also the case that

$$\forall \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \forall \sigma_2 \in \sigma_1 \quad \forall \sigma_3 \in \sigma_2 \quad F^*(\sigma_3)$$

is false since taking  $\sigma_1 = \{\{1\}, \{0, 1\}, \{0, 1, -1\}\}$ ,  $\sigma_2 = \{0, 1, -1\}$ ,  $\sigma_3 = -1$ , the value of  $F^*(\sigma_3)$  is false. □

**Proof of Proposition 11.14:** The proof is by induction on the number  $k$  of quantifiers, starting from the one outside.

Since  $(\forall X) \Phi$  is equivalent to  $\neg(\exists X) \neg\Phi$ , we can suppose without loss of generality that  $\text{Qu}_1$  is  $\exists$ .

The claim is certainly true when there is only one existential quantifier, by definition of  $\text{sign}(\mathcal{P})$ .

Suppose that

$$(\exists X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(X_1, \dots, X_k),$$



is true, and choose  $a \in \mathbb{R}$  such that

$$(\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(a, \dots, X_k)$$

is true. Note that, if  $\mathcal{P}_a$  is the set of polynomials obtained by substituting  $a \in \mathbb{R}$  to  $X_1$  in  $\mathcal{P}$ ,

$$\text{CSIGN}_1(\mathcal{P})(a) = \text{CSIGN}(\mathcal{P}_a).$$

By induction hypothesis,

$$(\text{Qu}_2 \sigma_2 \in \text{CSIGN}(\mathcal{P}_a)) \dots (\text{Qu}_k \sigma_k \in \sigma_{k-1}) F^*(\sigma_k) \text{ is true.}$$

is true. So, taking  $\sigma_1 = \text{CSIGN}(\mathcal{P}_a) = \text{CSIGN}(\mathcal{P})(a) \in \text{CSIGN}(\mathcal{P})$ ,

$$\exists \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} F^*(\sigma_k) \text{ is true.}$$

Conversely suppose

$$\exists \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} \quad F^*(\sigma_k)$$

is true and choose  $\sigma_1 \in \text{CSIGN}(\mathcal{P})$  such that

$$\text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} \quad F^*(\sigma_k)$$

is true. By definition of  $\text{CSIGN}(\mathcal{P})$ ,  $\sigma_1 = \text{CSIGN}(\mathcal{P})(a)$  for some  $a \in \mathbb{R}$ , and hence

$$\text{Qu}_2 \sigma_2 \in \text{CSIGN}(\mathcal{P}_a) \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} \quad F^*(\sigma_k)$$

is true. By induction hypothesis,

$$(\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(a, \dots, X_k)$$

is true. Thus

$$(\exists X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(X_1, \dots, X_k)$$

is true. □

Before giving a description of the Cylindrical Decision Algorithm, we explain how it works on the following example.

*Example 11.16.* We continue Example 11.10 to illustrate Proposition 11.14. We had determined

$$\text{CSIGN}(\mathcal{P}) = \{ \{ \{ \{ 1 \}, \{ 0, 1 \}, \{ 0, 1, -1 \} \}, \{ \{ 1 \}, \{ 0, 1 \} \}, \{ \{ 1 \} \} \}.$$

The formula

$$(\exists X_1) (\forall X_2) (\forall X_3) X_1^2 + X_2^2 + X_3^2 - 1 > 0$$

is certainly true since

$$\exists \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \forall \sigma_2 \in \sigma_1 \quad \forall \sigma_3 \in \sigma_2 \quad \sigma_3(P) = 1:$$

take  $\sigma_1 = \{\{1\}\}$ . It is also the case that the formula

$$(\forall X_1) (\exists X_2) (\exists X_3) X_1^2 + X_2^2 + X_3^2 - 1 > 0$$

is true since

$$\forall \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \exists \sigma_2 \in \sigma_1 \quad \exists \sigma_3 \in \sigma_2 \quad \sigma_3(P) = 1.$$

The formula

$$(\forall X_1) (\exists X_2) (\exists X_3) X_1^2 + X_2^2 + X_3^2 - 1 = 0$$

is false since it is not the case that

$$\forall \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \exists \sigma_2 \in \sigma_1 \quad \exists \sigma_3 \in \sigma_2 \quad \sigma_3(P) = 0:$$

take  $\sigma_1 = \{\{1\}\}$  to obtain a counter-example. It is also easy to check that the formula

$$(\exists X_1) (\forall X_2) (\exists X_3) X_1^2 + X_2^2 + X_3^2 - 1 = 0$$

is false since it is not the case that

$$\exists \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \forall \sigma_2 \in \sigma_1 \quad \exists \sigma_3 \in \sigma_2 \quad \sigma_3(P) = 0. \quad \square$$

We are ready for the Decision Algorithm using cylindrical decomposition. We consider a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , where  $D$  is an ordered integral domain.

*Algorithm 11.14. [Cylindrical Decision]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a  $\mathcal{P}$ -sentence

$$\Phi = (\text{Qu}_1 X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(X_1, \dots, X_k),$$

where  $F(X_1, \dots, X_k)$  is quantifier free,  $\text{Qu}_i \in \{\exists, \forall\}$ .

- **Output:** True if  $\Phi$  is true and False otherwise.
- **Procedure:**
  - Run Algorithm 11.2 (Cylindrical Decomposition) with input  $X_1, \dots, X_k$  and  $\mathcal{P}$  using Algorithm 11.13 for the Lifting Phase.
  - Extract  $\text{CSIGN}(\mathcal{P})$  from the set of cylindrical sample points and the signs of the polynomials in  $\mathcal{P}$  on the cells of  $\mathbb{R}^k$  using Proposition 11.9.
  - Trying all possibilities, decide whether

$$\text{Qu}_1 \sigma_1 \in \text{CSIGN}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} \quad F^*(\sigma_k) = \text{True},$$

which is clearly a finite verification.

**Proof of correctness:** Follows from Proposition 11.14. Note that the two first steps of the computation depend only on  $\mathcal{P}$  and not on  $\Phi$ . As noted before  $\text{CSIGN}(\mathcal{P})$  allows us to decide the truth or falsity of every  $\mathcal{P}$ -sentence.  $\square$

**Exercise 11.7.** Prove that the complexity of Algorithm 11.14 (Cylindrical Decision)  $(sd)^{O(1)^k}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{P}$ .

### 11.3 Quantifier Elimination

We start by explaining that the set of points of  $\mathbb{R}^\ell$  at which a  $\mathcal{P}$ -formula  $\Phi$  with free variables  $Y_1, \dots, Y_\ell$  is true, is a union of cells in  $\mathbb{R}^\ell$  of a cylindrical decomposition adapted to  $\mathcal{P}$ .

Indeed, let  $\mathcal{P} \subset \mathbb{R}[Y_1, \dots, Y_\ell, X_1, \dots, X_k]$  and let  $\mathcal{S}_1, \dots, \mathcal{S}_{\ell+k}$  a cylindrical decomposition of  $\mathbb{R}^{k+\ell}$  adapted to  $\mathcal{P}$ . Let  $S \in \mathcal{S}_i$ . We denote  $\text{CSIGN}_i(\mathcal{P})(y)$  for  $y \in S$  by  $\text{CSIGN}_i(\mathcal{P})(S)$ , using Proposition 11.11.

Let  $\Phi(Y) = (\text{Qu}_1 X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(Y_1, \dots, Y_\ell, X_1, \dots, X_k)$ , where  $F(Y_1, \dots, Y_\ell, X_1, \dots, X_k)$  is quantifier free,  $\text{Qu}_i \in \{\exists, \forall\}$ , be a  $\mathcal{P}$ -formula. Let  $\mathcal{L}$  be the union of cells  $S$  of  $\mathcal{S}_\ell$  such that

$$\text{Qu}_1 \sigma_1 \in \text{CSIGN}_\ell(\mathcal{P})(S) \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} F^*(\sigma_k) = \text{True}.$$

Then  $\text{Reali}(\Phi, \mathbb{R}^\ell) = \{y \in \mathbb{R}^\ell \mid \Phi(y)\} = \mathcal{L}$ . So we are not far from quantifier elimination.

However, a union of cells of a cylindrical decomposition in  $\mathbb{R}^\ell$  is not necessarily the realization of a  $C_{\leq \ell}(\mathcal{P})$ -quantifier free formulas, where  $C_{\leq \ell}(\mathcal{P}) = \bigcup_{i \leq \ell} C_i(\mathcal{P})$ . So a cylindrical decomposition does not always provide a  $C_{\leq \ell}(\mathcal{P})$ -quantifier free formula equivalent to  $\Phi$ . We give an example of this situation:

*Example 11.17.* Continuing Example 11.1 b), we consider  $\mathcal{P} = \{P, Q\}$  with  $P = X_2^2 - X_1(X_1 + 1)(X_1 - 2)$  and  $Q = X_2^2 - (X_1 + 2)(X_1 - 1)(X_1 - 3)$ .

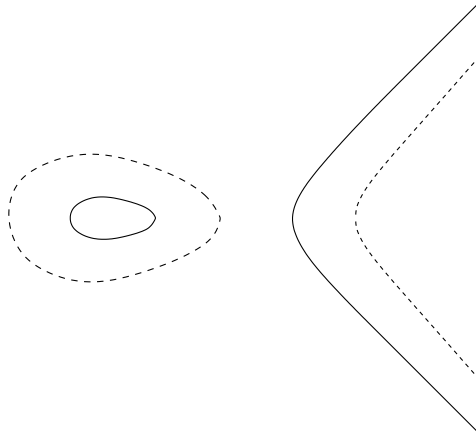
We have seen in Example 11.1 b) that

$$C_1(\mathcal{P}) = \{A, B, C\},$$

with

$$\begin{aligned} A(X_1) &= \text{sRes}_0(P, \partial P / \partial X_2) \\ &= 4 X_1 (X_1 + 1) (X_1 - 2), \\ B(X_1) &= \text{sRes}_0(Q, \partial Q / \partial X_2) \\ &= 4 (X_1 + 2) (X_1 - 1) (X_1 - 3), \\ C(X_1) &= \text{sRes}_0(P, Q) \\ &= (-X_1^2 - 3 X_1 + 6)^2. \end{aligned}$$

The zero sets of  $P$  and  $Q$  in  $\mathbb{R}^2$  are two cubic curves with no intersection (see Figure 11.2).



**Fig. 11.2.**  $\text{Zer}(P, \mathbb{R}^2)$  (solid) and  $\text{Zer}(Q, \mathbb{R}^2)$  (dotted)

This can be checked algebraically. The roots of  $(-X_1^2 - 3X_1 + 6)^2$ , which is the resultant of  $P$  and  $Q$ , are  $a = -3/2 + (1/2)\sqrt{33}$  and  $b = 3/2 - (1/2)\sqrt{33}$ . Substituting these values in  $P$  and  $Q$  gives polynomials of degree 2 without real roots.

The only subset of  $\mathbb{R}$  defined by sign conditions on  $C_1(\mathcal{P})$  are

$$\begin{aligned}
 \{-1, 0\} &= \{x \in \mathbb{R} \mid A(x) = 0 \wedge B(x) > 0 \wedge C(x) > 0\}, \\
 (-1, 0) \cup (3, +\infty) &= \{x \in \mathbb{R} \mid A(x) > 0 \wedge B(x) > 0 \wedge C(x) > 0\}, \\
 (-2, -1) \cup (0, 1) &= \{x \in \mathbb{R} \mid A(x) < 0 \wedge B(x) > 0 \wedge C(x) > 0\}, \\
 \{3\} &= \{x \in \mathbb{R} \mid A(x) > 0 \wedge B(x) = 0 \wedge C(x) > 0\}, \\
 \{-2, 1\} &= \{x \in \mathbb{R} \mid A(x) < 0 \wedge B(x) = 0 \wedge C(x) > 0\}, \\
 \{2\} &= \{x \in \mathbb{R} \mid A(x) = 0 \wedge B(x) < 0 \wedge C(x) > 0\}, \\
 (2, 3) &= \{x \in \mathbb{R} \mid A(x) > 0 \wedge B(x) < 0 \wedge C(x) > 0\}, \\
 (-\infty, -2) \cup (1, 2) \setminus \{a, b\} &= \{x \in \mathbb{R} \mid A(x) < 0 \wedge B(x) < 0 \wedge C(x) > 0\}, \\
 \{a, b\} &= \{x \in \mathbb{R} \mid A(x) < 0 \wedge B(x) < 0 \wedge C(x) = 0\}.
 \end{aligned}$$

The set  $\{x \in \mathbb{R} \mid \exists y \in \mathbb{R} P(x, y) < 0 \wedge Q(x, y) > 0\} = (2, +\infty)$  is the union of semi-algebraically connected components of semi-algebraic sets defined by sign conditions on  $C_1(\mathcal{P})$ , but is not defined by any  $C_1(\mathcal{P})$ -quantifier free formula. There are  $\mathcal{P}$ -formulas whose realization set cannot be described by  $C_1(\mathcal{P})$ -quantifier free formulas.  $\square$

Fortunately, closing the set of polynomials under differentiation before each application of elimination of a variable provides an extended cylindrical family whose realization of sign conditions are the cells of a cylindrical decomposition. This has been already proved in Theorem 5.34,

We denote by  $\overline{C}_k(\mathcal{P})$  the set of polynomials in  $\mathcal{P}$  and all their derivatives with respect to  $X_k$ , and by  $\overline{C}_i(\mathcal{P})$  the set obtained by adding to the polynomials in  $\text{Elim}_{X_{i+1}}(\overline{C}_{i+1}(\mathcal{P}))$ , all their derivatives with respect to  $X_i$ , so that  $\overline{C}_i(\mathcal{P}) \subset \mathbb{R}[X_1, \dots, X_i]$ . According to Theorem 5.34, the realization of sign conditions on  $\overline{C}_{\leq i}(\mathcal{P}) = \bigcup_{j \leq i} \overline{C}_j(\mathcal{P})$  are the sets of a cylindrical decomposition of  $\mathbb{R}^i$  and the realization of sign conditions on  $\overline{C}(\mathcal{P}) = \overline{C}_{\leq k}(\mathcal{P})$  are the sets of a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .

*Algorithm 11.15. [Improved Cylindrical Decomposition]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** an ordered list of variables  $X_1, \dots, X_k$ , a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** the finite set of polynomials  $\overline{C}_i(\mathcal{P}) \subset D[X_1, \dots, X_i]$  and the realizable sign conditions on  $\overline{C}_{\leq i}(\mathcal{P})$  for every  $i = k, \dots, 1$ . The non-empty realizations of sign conditions on  $\overline{C}_{\leq i}(\mathcal{P})$ ,  $i = 1, \dots, k$  constitute a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ .
- **Procedure:**
  - Add to the elements of  $\mathcal{P}$  all their derivatives with respect to  $X_k$ , which defines  $\overline{C}_k(\mathcal{P})$ ,
  - Elimination phase: Compute  $\overline{C}_i(\mathcal{P})$  for  $i = k - 1, \dots, 1$ , using Algorithm 11.1 (Elimination) and adding the derivatives with respect to  $X_i$ .
  - Lifting phase:
    - Compute the sample points of the cells in  $\mathcal{S}_1$  by characterizing the roots of  $\overline{C}_1(\mathcal{P})$  and choosing a point in each interval they determine, using Algorithm 11.13 (Lifting Phase) or Algorithm 11.7 (Real Lifting Phase).
    - For every  $i = 2, \dots, k$ , compute the sample points of the cells of  $\mathcal{S}_i$  from the sample points of the cells in  $\mathcal{S}_{i-1}$  as follows: Compute for every sample point  $x$  of the cells in  $\mathcal{S}_{i-1}$ , the list  $L$  of non-zero polynomials  $P_i(x, X_i)$  with  $P_i \in \overline{C}_i(\mathcal{P})$ , using Algorithm 11.4 (Real Recursive Sign at a Point) or Algorithm 11.8 (Recursive Sign Determination). Characterize the roots of  $L$  and choose a point in each interval they determine using Algorithm 11.13 (Lifting Phase) or Algorithm 11.7 (Real Lifting Phase).
  - Output the sample points of the cells with the sign condition on  $\overline{C}_{\leq i}(\mathcal{P})$  valid at the sample point of each cell of  $\mathbb{R}^i$ .

**Proof of correctness:** Follows from Theorem 5.34. Note that the realization of sign conditions on  $\overline{C}_{\leq i}(\mathcal{P})$  are semi-algebraically connected subsets of  $\mathbb{R}^i$ .  $\square$

**Exercise 11.8.** Prove that the complexity of Algorithm 11.15 (Improved Cylindrical Decision) is  $(sd)^{O(1)k}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{P}$ .

We are going to see that the Improved Cylindrical Decomposition Algorithm with input  $\mathcal{P}$  makes it possible to eliminate quantifiers of any  $\mathcal{P}$ -formula. We need the following notation:

For every non-empty sign condition  $\sigma$  on  $\overline{C}_{i \leq \ell}(\mathcal{P})$ ,  $\text{CSIGN}_{\ell}(\mathcal{P})(x)$  is constant as  $x$  varies in the realization of  $\sigma$ , by Proposition 11.11, and is denoted by  $\text{CSIGN}_{\ell}(\mathcal{P})(\sigma)$ .

*Algorithm 11.16. [Cylindrical Quantifier Elimination]*

- **Structure:** an integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a finite set  $\mathcal{P} \subset D[Y_1, \dots, Y_{\ell}][X_1, \dots, X_k]$ , a  $\mathcal{P}$ -formula

$$\Phi(Y) = (\text{Qu}_1 X_1) (\text{Qu}_2 X_2) \dots (\text{Qu}_k X_k) F(Y_1, \dots, Y_{\ell}, X_1, \dots, X_k).$$

where  $F(Y_1, \dots, Y_{\ell}, X_1, \dots, X_k)$  is quantifier free,  $\text{Qu}_i \in \{\exists, \forall\}$ , with free variables  $Y = Y_1, \dots, Y_{\ell}$ .

- **Output:** a quantifier free formula  $\Psi(Y)$  equivalent to  $\Phi(Y)$ .
- **Procedure:**
  - Run Algorithm 11.15 (Improved Cylindrical Decomposition) with input  $Y_1, \dots, Y_{\ell}, X_1, \dots, X_k$  and  $\mathcal{P}$ .
  - For every non-empty sign condition  $\sigma$  on  $\overline{C}_{\leq \ell}(\mathcal{P})$ , extract  $\text{CSIGN}_{\ell}(\mathcal{P})(\sigma)$  from the sample points of the cells and the signs of the polynomials in  $\mathcal{P}$  on the cells of  $R^k$  using Proposition 11.9.
  - Make the list  $\mathcal{L}$  of the non-empty sign condition  $\sigma$  on  $\overline{C}_{\leq \ell}(\mathcal{P})$  for which
 
$$\text{Qu}_1 \sigma_1 \in \text{CSIGN}_{\ell}(\mathcal{P})(\sigma) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_k \sigma_k \in \sigma_{k-1} \quad F^*(\sigma_k) = \text{True}.$$
  - Output

$$\Psi(Y) = \bigvee_{\sigma \in \mathcal{L}} \bigwedge_{P \in C_{\leq \ell}(\mathcal{P})} \text{sign}(P(Y_1, \dots, Y_{\ell})) = \sigma(P).$$

**Proof of correctness:** Follows from Theorem 5.34. □

**Exercise 11.9.** Prove that the complexity of Algorithm 11.16 (Cylindrical Quantifier Elimination) is  $(s d)^{O(1)^k}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{P}$ . When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , prove the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(1)^{k-1}}$ .

## 11.4 Lower Bound for Quantifier Elimination

In this section, we prove that a doubly exponential complexity for the quantifier elimination problem is unavoidable. We first need a notion for the size of a formula, which we define as follows.

We define the size of atomic formulas  $P > 0, P <, P = 0$  to be the number of coefficients needed to write the polynomial  $P$  in the dense form. Thus, if  $P \in \mathbb{R}[X_1, \dots, X_k]$  and  $\deg(P) = d$ ,

$$\text{size}(P > 0) = \text{size}(P = 0) = \text{size}(P < 0) := \binom{d+k}{k}$$

by Lemma 8.5.

Next we define inductively, for formulas  $\phi_1, \phi_2$ ,

$$\begin{aligned} \text{size}(\phi_1 \vee \phi_2) = \text{size}(\phi_1 \wedge \phi_2) &:= \text{size}(\phi_1) + \text{size}(\phi_2) + 1, \\ \text{size}(\neg \phi_1) &:= \text{size}(\phi_1) + 1, \\ \text{size}(\exists X \phi_1) = \text{size}(\forall X \phi_1) &:= \text{size}(\phi_1) + 2. \end{aligned}$$

We prove the following theorem.

**Theorem 11.18.** *There exist natural numbers  $c, c' > 0$  and a sequence of quantified formulas  $\phi_n(X, Y)$  such that  $\text{size}(\phi_n) \leq c'n$  and any quantifier-free formula equivalent to  $\phi_n$  must have size at least  $2^{c+2^{n-3}}$ .*

**Proof:** We construct the formulas  $\phi_n(X, Y)$  as follows. It is useful to consider  $Z = X + iY$  as a complex variable. We now define a predicate,  $\psi_n(W, Z)$  such that  $\psi_n(W, Z)$  holds if and only if  $W = Z^{2^{2^n}}$ . Here, both  $W$  and  $Z$  should be thought of as complex variables. The predicate  $\psi_n$  is defined recursively as follows:

$$\begin{aligned} \psi_0(W, Z) &:= (W - Z^2 = 0), \\ \psi_n(W, Z) &:= (\exists U)(\forall A \forall B)((A = W \wedge B = U) \vee (A = U \wedge B = Z)) \\ &\quad \Rightarrow \psi_{n-1}(A, B). \end{aligned} \tag{11.2}$$

It is easy to check that formula (11.2) is equivalent to formula,

$$(\exists U)\psi_{n-1}(W, U) \wedge \psi_{n-1}(U, Z), \tag{11.3}$$

which is clearly equivalent to  $W = Z^{2^{2^n}}$ . Moreover the recursion in formula (11.2) implies that  $\text{size}(\psi_n(W, Z)) \leq c_1 n$ , where  $c_1$  is a natural number.

We now define  $\phi_n(X, Y)$  to be the formula obtained by specializing  $W$  to 1 in the formula  $\psi_n$ , as well as writing the various complex variables appearing in the formula in terms of their real and imaginary parts. It is easy to check that  $\text{size}(\phi_n) \leq c'n$  where  $c'$  is a natural number,  $c' \geq c_1$ .

Now, let  $\theta_n(X, Y)$  be a quantifier-free formula equivalent to  $\phi_n(X, Y)$ . Let  $\mathcal{P}_n = \{P_1, \dots, P_s\}$  denote the set of polynomials appearing in  $\theta_n$  and let  $\deg(P_i) = d_i$ . From the definition of the size of a formula we have,

$$\text{size}(\theta_n) \geq \sum_{i=1}^s d_i.$$

Clearly, the set  $S \subset \mathbb{R}^2$  defined by  $\theta_n$  has  $2^{2^n}$  isolated points (corresponding to the different  $2^{2^n}$ -th complex roots of unity). But  $S$  is a  $\mathcal{P}_n$ -semi-algebraic set. By Theorem 7.50, there exists a natural number  $C$  such that

$$\sum_{0 \leq i \leq 2} b_i(S) \leq C (s d)^4,$$

where  $d = \sum_{1 \leq i \leq s} d_i$  is an upper bound on the degrees of the polynomials in  $\mathcal{P}_n$ . Moreover,  $s \leq d$ . Thus, we have that,  $b_0(S) = 2^{2^n} \leq C (s d)^4 \leq C d^8$ . Hence,  $\text{size}(\theta_n) \geq d \geq 2^{c+2^{n-3}}$ .  $\square$

Notice however in the above proof the number of quantifier alternations in the formulas  $\phi_n$  is linear in  $n$ . Later in Chapter 14, we will develop an algorithm for performing quantifier elimination whose complexity is doubly exponential in the number of quantifier alternations, but whose complexity is only singly exponential if we fix the number of quantifier alternations allowed in the input formula.

## 11.5 Computation of Stratifying Families

When we want to decide the truth of a sentence, or eliminate quantifiers from a formula, the variables provided in the input play a special role in the problem considered and cannot be changed. However when we are only interested in computing the topology of a set, we are free to perform linear changes of coordinates.

We indicate now how to compute a cell stratification of  $\mathbb{R}^k$  adapted to a finite set of polynomials  $\mathcal{P}$  using Section 5.5 of Chapter 5.

### *Algorithm 11.17.* [Stratifying Cylindrical Decomposition]

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$  of  $s$  polynomials of degree bounded by  $d$ .
- **Output:** a cell stratification of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ . More precisely, a linear automorphism  $u$  such that the finite set of polynomials

$$\overline{C}_i(u(\mathcal{P})) \subset D[X_1, \dots, X_i]$$

are quasi-monic with respect to  $X_i$  and the realizable sign conditions on  $\overline{C}_{\leq i}(u(\mathcal{P}))$  for every  $i = 1, \dots, k$ . The families  $\mathcal{S}_i$ , for  $i = 1, \dots, k$ , consisting of all  $\text{Reali}(\sigma)$  with  $\sigma$  a realizable sign conditions on  $\overline{C}_{\leq i}(u(\mathcal{P}))$  constitute a cell stratification of  $\mathbb{R}^k$  adapted to  $u(\mathcal{P})$ .



- **Procedure:**

- Try successively

$$a_k = (a_{k,1}, \dots, a_{k,k-1}) \in \{0, \dots, s d\}^{k-1}$$

and choose one such that after the linear change of variables  $u_k$  associating to

$$X_1, X_2, \dots, X_{k-1}, X_k,$$

$$X_1 + a_{k,1} X_k, X_2 + a_{k,2} X_k, \dots, X_{k-1} + a_{k,k-1} X_k, X_k$$

the polynomials in  $u_k(\mathcal{P})$  are monic in  $X_k$ .

- Add to the elements of  $u_k(\mathcal{P})$  all their derivatives with respect to  $X_k$ , which defines  $\overline{C}_k(u_k(\mathcal{P}))$ .
- **Elimination phase:**

For  $i = k - 1, \dots, 1$ , denote by  $d_i$  and  $s_i$  a bound on the degree and number of the polynomials in  $\overline{C}_{i+1}(u_{i+1}(\mathcal{P}))$  and choose

$$a_i = (a_{i,1}, \dots, a_{i,i-1}) \in \{0, \dots, s_i d_i\}^{i-1}$$

such that after the linear change of variables  $v_i$  associating to

$$X_1, X_2, \dots, X_{i-1}, X_i, \dots, X_k,$$

$$X_1 + a_{i,1} X_i, X_2 + a_{i,2} X_i, \dots, X_{i-1} + a_{i,i-1} X_i, X_i, \dots, X_k$$

the polynomials in  $\overline{C}_{i+1}(u_i(\mathcal{P}))$  are monic in  $X_i$ , with  $u_i = v_i \circ u_{i+1}$ . Compute  $\overline{C}_i(u_i(\mathcal{P}))$  for  $i = k - 1, \dots, 1$ , using Algorithm 11.1 (Elimination) and adding the derivatives with respect to  $X_i$ . Define  $u = u_1 = v_1 \circ \dots \circ v_k$ .

- **Lifting phase:**
  - Compute the sample points of the cells in  $\mathcal{S}_1$ , by characterizing the roots of  $C_1(u(\mathcal{P}))$  and choosing a point in each interval they determine using Algorithm 11.13 (Lifting Phase) or Algorithm 11.7 (Real Lifting Phase).
  - For every  $i = 2, \dots, k$ , compute the sample points of the cells of  $\mathcal{S}_i$  from the sample points of the cells in  $\mathcal{S}_{i-1}$ , as follows: Compute, for every sample point  $x$  of a cell in  $\mathcal{S}_{i-1}$ , the list  $L$  of non-zero polynomials  $Q \in \overline{C}_i(u(\mathcal{P}))$  using Algorithm 11.4 (Real Recursive Sign at a Point) or Algorithm 11.8 (Recursive Sign Determination). Characterize the roots of  $L$  and chose a point in each interval they determine using Algorithm 11.13 (Lifting Phase) or Algorithm 11.7 (Real Lifting Phase).
  - Output the sample points of the cells with the sign condition on  $\overline{C}_{\leq i}(u(\mathcal{P}))$  valid at the sample point of each cell of  $\mathbb{R}^j$ ,  $j \leq i$ .

**Proof of correctness:** Follows from Lemma 4.73, Lemma 4.74, Section 5.5 of Chapter 5, Lemma 4.74, and the correctness of Algorithm 11.1 (Elimination).  $\square$

**Exercise 11.10.** Prove that the complexity of Algorithm 11.17 (Stratifying Cylindrical Decision) is  $(sd)^{O(1)k}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{P}$ . When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , prove the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(1)k-1}$ .

Using cylindrical decomposition, it is possible to design algorithms for computing various topological informations.

*Remark 11.19.*

a) It is possible to compute a semi-algebraic description of the semi-algebraically connected components of a  $\mathcal{P}$ -semi-algebraic set with complexity  $(sd)^{O(1)k}$  as follows, using the proofs of Theorem 5.21 and Theorem 5.38 and the complexity analysis of Algorithm 11.17 (Stratifying Cylindrical Decomposition) (Exercise 11.10). Hint: Consider a  $\mathcal{P}$ -semi-algebraic set, compute a stratifying family adapted to  $\mathcal{P}$  and a description of the cells associated to the corresponding Stratifying Cylindrical Decomposition. Determine the adjacency relations between cells from the list of non-empty sign conditions, using Proposition 5.39 (Generalized Thom's Lemma). Use these adjacencies to describe the connected components of  $S$ . We do not give details since we are going to see in Chapters 15 and 16 much better algorithms for the description of connected components of semi-algebraic sets.

b) A triangulation of a closed and bounded semi-algebraic set, as well as its homology groups, can be computed with complexity  $(sd)^{O(1)k}$ , using the proof of Theorem 5.43, the definition of the homology groups of bounded and closed semi-algebraic sets in Chapter 6, and the complexity analysis of Algorithm 11.17 (Stratifying Cylindrical Decomposition) (Exercise 11.10). We do not give details either, even though we currently have no better algorithm for computing the homology groups.

c) A bound on the finite number of topological types of algebraic sets defined by polynomials of degree  $d$  follows from Algorithm 11.17 (Stratifying Cylindrical Decomposition). The bound is polynomial in  $d$  and doubly exponential in the number  $\binom{d+k}{k}$  of monomials of degree  $d$  in  $k$  variables. Again, we do not give details for this quantitative version of Theorem 5.47.  $\square$

## 11.6 Topology of Curves

In this section,  $D$  is an ordered integral domain contained in a real closed field  $\mathbb{R}$  and  $C = \mathbb{R}[i]$ .

The simplest situation where the Cylindrical Decomposition Algorithm can be performed is the case of one single non-zero polynomial bivariate polynomial  $P(X, Y) \in D[X, Y]$ .

The zero set of this polynomial is an algebraic subset  $\text{Zer}(P, \mathbb{R}^2) \subset \mathbb{R}^2$  contained in the plane  $\mathbb{R}^2$  and distinct from  $\mathbb{R}^2$ .

Since an algebraic set is semi-algebraic, there are three possible cases:

- $\text{Zer}(P, \mathbb{R}^2)$  is an algebraic set of dimension 1, which is called a **curve**.
- $\text{Zer}(P, \mathbb{R}^2)$  is an algebraic set of dimension 0, i.e. a finite set of points.
- $\text{Zer}(P, \mathbb{R}^2)$  is empty.

Typical examples of these three situations are the unit circle defined by  $X^2 + Y^2 - 1 = 0$ , the origin defined by  $X^2 + Y^2 = 0$ , and the empty set defined by  $X^2 + Y^2 + 1 = 0$ .

Let us consider a polynomial  $P(X, Y)$  in two variables, monic and separable, of degree  $d$  as a polynomial in  $Y$ . By separable we mean in this section that any gcd of  $P$  and  $\partial P / \partial Y$  is an element of  $\mathbb{R}(X)$ . This is not a big loss of generality since it is always possible to make a polynomial monic by a change of variables of the form  $X + aY$ ,  $a \in \mathbb{Z}$ , by Lemma 4.73. Replacing  $P$  by a separable polynomial can be done using Algorithm 10.1 (Gcd and gcd-free part), and does not modify the zero set. The zero set of this polynomial is an algebraic set  $\text{Zer}(P, \mathbb{R}^2)$  contained in the plane.

The cylindrifying family of polynomials associated to  $P$  consists of the subdiscriminants of  $P$  with respect to  $Y$  (see page 102 and Proposition 4.27), which are polynomials in the variable  $X$ , denoted by  $\text{sDisc}_j(X)$  for  $j$  from 0 to  $d - 1$  (since  $P$  is monic in  $Y$ ,  $\text{sDisc}_d$  and  $\text{sDisc}_{d-1}$  are constant). Denote by  $\text{sDisc}$  the list  $\text{sDisc}_d, \dots, \text{sDisc}_0$ . Note that  $\text{Disc}(X) = \text{sDisc}_0(X)$  is the discriminant, and is not identically 0 since  $P$  is separable (see Equation (4.4) and Corollary 4.2). On intervals between the roots of  $\text{Disc}_0(X)$ , the number of roots of  $P(x, Y)$  in  $\mathbb{R}$  is fixed (this is a special case of Theorem 5.16). The number of roots of  $P(x, Y)$  can be determined by the signs of the other signed subresultant coefficients and is equal to  $\text{PmV}(\text{sDisc})$  according to Theorem 4.33. Note that on an interval between two roots of  $\text{Disc}(X)$  in  $\mathbb{R}$ , the signs of the subdiscriminants may change but  $\text{PmV}(\text{sDisc})$  is fixed.

A cylindrical decomposition of the plane adapted to  $P$  can thus be obtained as follows: the cells of  $\mathbb{R}$  are the roots of  $\text{Disc}(X)$  in  $\mathbb{R}$  and the intervals they determine. Above each root of  $\text{Disc}(X)$  in  $\mathbb{R}$ ,  $\text{Zer}(P, \mathbb{R}^2)$  contains a finite number of points. These points and intervals between them are cells of  $\mathbb{R}^2$ . Above each interval determined by the roots of  $\text{Disc}(X)$ , there are 1 dimensional cells, called **curve segments**, which are graphs of functions from the interval to  $\mathbb{R}$ , and 2 dimensional cells which are bands between these graphs. Finally,  $\text{Zer}(P, \mathbb{R}^2)$  is the union of the points and curve segments of this cylindrical decomposition and consists of a finite number of points projecting on the roots of  $\text{Disc}(X)$  and a finite number of curve segments homeomorphic to segments of  $\mathbb{R}$  and projecting on intervals determined by the roots of  $\text{Disc}(X)$ . If above every interval determined by the roots of  $\text{Disc}(X)$  there are no curve segments,  $\text{Zer}(P, \mathbb{R}^2)$  is at most a finite number of points, or empty. If moreover above every root of  $\text{Disc}(X)$  there are no points  $\text{Zer}(P, \mathbb{R}^2)$  is empty.

The purpose of the algorithm we are going to present is to compute exactly the topology of the curve  $\text{Zer}(P, \mathbb{R}^2)$ , i.e. to determine a planar graph homeomorphic to  $\text{Zer}(P, \mathbb{R}^2)$ . After performing, if necessary, a linear change of coordinates, this will be done by indicating adjacencies between points of the curve  $\text{Zer}(P, \mathbb{R}^2)$  above roots of  $\text{Disc}(X)$  and curve segments on intervals between these roots. In this study, the notions of critical points of the projection and of curves in generic position will be useful.

The **critical points** of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  to the  $X$ -axis are the points  $(x, y) \in \text{Zer}(P, \mathbb{C}^2)$  such that  $y$  is a multiple root of  $P(x, Y)$ .

The critical points of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  to the  $X$ -axis are of two kinds

- **singular points** of  $\text{Zer}(P, \mathbb{C}^2)$ , i.e. points of  $\text{Zer}(P, \mathbb{C}^2)$  where

$$\partial P / \partial X(x, y) = \partial P / \partial Y(x, y) = 0,$$

- **ordinary critical points** points where the tangent to  $\text{Zer}(P, \mathbb{C}^2)$  is well defined and parallel to the  $Y$ -axis, i.e. points  $(x, y) \in \mathbb{C}^2$  where

$$\partial P / \partial Y(x, y) = 0, \partial P / \partial X(x, y) \neq 0.$$

In both cases, the first coordinate of a critical point of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  to the  $X$ -axis is a root of the discriminant  $\text{Disc}(X)$  of  $P$  considered as a polynomial in  $Y$ .

Computing the topology will be particularly easy for a curve in generic position, which is the notion we define now. Indeed, the critical points of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  on the  $X_1$ -axis are easy to characterize in this case.

Let  $P$  be polynomial of degree  $d$  in  $\mathbb{R}[X, Y]$  that is separable. The set  $\text{Zer}(P, \mathbb{C}^2)$  is in **generic position** if the following two conditions are satisfied:

- $\deg(P) = \deg_Y(P)$ ,
- for every  $x \in \mathbb{C}$ ,  $\gcd(P(x, Y), \partial P / \partial Y(x, Y))$  is either a constant or a polynomial of degree  $j$  with exactly one root of multiplicity  $j$ . In other words, there is at most one critical point  $(x, y)$  of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  to the  $X_1$ -axis above any  $x \in \mathbb{C}$ .

Note that above an element  $x_1$  of  $\mathbb{R}$  which is a root of  $\text{Disc}(X)$ , the unique root of  $\gcd(P(x, Y), \partial P / \partial Y(x, Y))$  is necessarily in  $\mathbb{R}$ . So there is exactly a critical point with coordinates in  $\mathbb{R}$  above each root of  $\text{Disc}(X)$  in  $\mathbb{R}$ .

*Example 11.20.* If  $P = (X^2 - Y + 1)(X^2 + Y^2 - 2Y)$ , the set  $\text{Zer}(P, \mathbb{C}^2)$  is not in generic position since there are two critical points of the projection on  $X$  above the point 0, namely  $(0, 0)$  and  $(0, 1)$ .  $\square$

The output of the algorithm computing the topology of a curve in generic position will be the following:

- the number  $r$  of roots of  $\text{Disc}(X)$  in  $\mathbb{R}$ . We denote these roots by  $x_1 < \dots < x_r$ , and by  $x_0 = -\infty, x_{r+1} = +\infty$ .

- the number  $m_i$  of roots of  $P(x, Y)$  in  $\mathbb{R}$  when  $x$  varies on  $(x_i, x_{i+1})$ , for  $i = 0, \dots, r$ .
- the number  $n_i$  of roots of  $P(x, Y)$  in  $\mathbb{R}$ . We denote these roots by  $y_{i,j}$ , for  $j = 1, \dots, n_i$ .
- a number  $c_i \leq n_i$  such that if  $C_i = (x_i, z_i)$  is the unique critical point of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  on the  $X$ -axis above  $x_i$ ,  $z_i = y_{i,c_i}$ .

More precisely, the output is

$$[m_0, [n_1, c_1], \dots, m_{r-1}, [n_r, c_r], m_r].$$

It is clear that a graph homeomorphic to  $\text{Zer}(P, \mathbb{R}^2) \subset \mathbb{R}^2$  can be drawn using the output since if  $m_i \geq n_i$  (resp.  $m_i \geq n_{i+1}$ ),  $C_i$  belong to the closure of  $m_i - n_i$  (resp.  $m_i - n_{i+1}$ ) curve segments above  $(x_i, x_{i+1})$  getting glued at  $C_i$  (resp.  $C_{i+1}$ ) and if  $m_i = n_i - 1$  (resp.  $m_i = n_{i+1} - 1$ ) the point  $C_i$  does not belong to the closure of a curve segment above  $(x_i, x_{i+1})$ .

The critical points of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  on the  $X$ -axis in the generic position case are easy to determine.

We denote

$$\begin{aligned} \text{sDiscP}_j(X, Y) &= \text{sResP}_j(P(X, Y), \partial P / \partial Y(X, Y)) \\ \text{sDisc}_j(X) &= \text{sRes}_j(P(X, Y), \partial P / \partial Y(X, Y)) \end{aligned}$$

( $P$  is considered as a polynomial in  $Y$ ).

**Proposition 11.21.** *Let  $P$  be a polynomial of degree  $d$  in  $\mathbb{R}[X, Y]$ , separable and in generic position. If  $(x, y)$  is a critical point of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  on the  $X$ -axis, and  $j$  is the multiplicity of  $y$  as a root of  $P(x, Y)$  (considered as a polynomial in  $X_2$ ), then*

$$\text{Disc}(x) = \text{sDisc}_1(x) = 0, \dots, \text{sDisc}_{j-1}(x) = 0, \text{sDisc}_j(x) \neq 0,$$

and

$$y = -\frac{1}{j} \frac{\text{sDisc}_{j,j-1}(x)}{\text{sDisc}_j(x)},$$

where  $\text{sDisc}_{j,j-1}(X)$  is the coefficient of  $Y^{j-1}$  in  $\text{sDiscP}_j(X, Y)$ .

**Proof:** Since  $(x, y) \in \mathbb{R}^2$  is a critical point of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  on the  $X$ -axis,  $j > 0$ . Since  $P$  is in generic position,  $\text{sDiscP}_j(x, Y)$  is a degree  $j$  polynomial with only one root,  $y$ , which implies that

$$\text{sDiscP}_j(x, Y) = \text{sDisc}_j(x)(Y - y)^j, \tag{11.4}$$

and

$$y = -\frac{1}{j} \frac{\text{sDisc}_{j,j-1}(x)}{\text{sDisc}_j(x)},$$

identifying coefficients of degree  $j - 1$  on both sides of (11.4). □

We can also count the number of points of the curve above the critical points using the following result.

**Proposition 11.22.** *Let  $P(X, Y) \in \mathbb{D}[X, Y]$  be a degree  $d$  separable polynomial such that  $\text{lcof}_Y(P) \in \mathbb{D}$  and  $x \in \mathbb{R}$ . Let*

$$R(X, Y, Z) = (Y - Z) \partial P / \partial Y(X, Y) - d P(X, Y).$$

We denote

$$T_j(X, Z) = \text{sRes}_j(P(X, Y), R(X, Y, Z))$$

(considered as polynomials in  $Y$ ) and  $T(X, Z)$  the list  $T_j(X, Z)$ ,  $j$  from 0 to  $d - 1$ . Let  $r = \text{PmV}(\text{sDisc}(x))$ .

a)

$$\#\{y \in \mathbb{R} \mid P(x, y) = 0\} = r$$

b) If  $P(x, z) \neq 0$ ,

$$\#\{y \in \mathbb{R} \mid P(x, y) = 0 \wedge z < y\} = (r + \text{PmV}(T(x, z))) / 2.$$

c) If  $P(x, z) = 0$ ,

$$\#\{y \in \mathbb{R} \mid P(x, y) = 0 \wedge z < y\} = (r + \text{PmV}(T(x, z)) - 1) / 2.$$

**Proof:** We denote

$$\begin{aligned} Z_x &= \{y \in \mathbb{R} \mid P(x, y) = 0\}, \\ \text{TaQ}(1, Z_x) &= \#(Z_x) \\ &= \#(\{y \in \mathbb{R} \mid P(x, y) = 0\}) \\ \text{TaQ}(Y - z, Z_x) &= \#(\{y \in \mathbb{R} \mid P(x, y) = 0 \wedge y > z\}) \\ &\quad - \#(\{y \in \mathbb{R} \mid P(x, y) = 0 \wedge y < z\}), \\ \text{TaQ}((Y - z)^2, Z_x) &= \#\{y \in \mathbb{R} \mid P(x, y) = 0 \wedge y \neq z\}. \end{aligned}$$

a) follows from Theorem 4.33. b) and c) follow from Equation (2.6) (see page 65) applied to  $P(x, Y)$  and  $Y - z$ :

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} c(Y = z, P(x, Y) = 0) \\ c(Y > z, P(x, Y) = 0) \\ c(Y < z, P(x, Y) = 0) \end{bmatrix} = \begin{bmatrix} \text{TaQ}(1, P(x, Y)) \\ \text{TaQ}(Y - z, P(x, Y)) \\ \text{TaQ}((Y - z)^2, P(x, Y)) \end{bmatrix}.$$

If  $P(x, z) \neq 0$ , we have  $c(Y - z = 0, P(x, Y) = 0) = 0$ . Thus

$$c(Y > z, P(x, Y) = 0) = (\text{TaQ}(1, P(x, Y)) + \text{TaQ}(Y - z, P(x, Y))) / 2.$$

If  $P(x, z) = 0$ , we have  $c(Y - z = 0, P(x, Y) = 0) = 1$ . Thus

$$c(Y > z, P(x, Y) = 0) = (\text{TaQ}(1, P(x, Y)) + \text{TaQ}(Y - z, P(x, Y)) - 1) / 2. \quad \square$$

Next we show that it is always possible to perform a linear change of variables such that in the new coordinates the curve is in generic position. The idea is to maximize the number of distinct roots of the discriminant, which are the  $X$ -coordinates of the critical points.

Let  $P$  be a separable polynomial in  $D[X, Y]$ ,  $U$  a new variable and  $Q(U, X, Y)$  the polynomial defined by:

$$Q(U, X, Y) = P(X + UY, Y).$$

If  $a \in \mathbb{Z}$ , let  $\text{Zer}(P_a, \mathbb{R}^2)$  denote the curve defined by the polynomial

$$P_a(X, Y) = Q(a, X, Y) = 0.$$

We denote

$$\begin{aligned} \text{sDiscP}_j(U, X, Y) &= \text{sResP}_j(Q(U, X, Y), \partial Q / \partial Y(U, X, Y)) \\ \text{sDisc}_j(U, X) &= \text{sRes}_j(Q(U, X, Y), \partial Q / \partial Y(U, X, Y)) \end{aligned}$$

Note that  $\text{Disc}(U, X) = \text{sDisc}_0(U, X)$  is the discriminant of  $Q(U, X, Y)$  with respect to the variable  $Y$ . We denote by  $\text{sDisc}(U, X)$  the list of the subdiscriminants  $\text{sDisc}_j(U, X)$ . Let  $\Delta(U)$  be the non-zero subdiscriminant of smallest possible index of  $\text{Disc}(U, X)$  with respect to the variable  $X$ .

**Proposition 11.23.** *Let  $a$  be an integer number such that*

$$\deg_{X_2}(Q(a, X, Y)) = \deg(Q(a, X, Y)) \quad \text{and} \quad \Delta(a) \neq 0.$$

*Then  $P_a$  is in generic position.*

**Proof:** Suppose that  $\deg_{X_2}(P_a(X, Y)) = \deg(P_a(X, Y))$  and  $P_a$  is not in generic position. Let  $\delta$  be a new variable and consider the field  $\mathbb{C}\langle\delta\rangle$  of algebraic Puiseux series in  $\delta$ . We are going to prove the following property (P): the number of distinct roots of  $\text{Disc}(a + \delta, X)$  in  $\mathbb{C}\langle\delta\rangle$ , which is the discriminant of  $P_{a+\delta}(X, Y)$  with respect to the variable  $Y$ , is bigger than the number of distinct roots of  $\text{Disc}(a, X)$  in  $\mathbb{C}$ , which is the discriminant of  $P_a(X, Y)$  with respect to the variable  $Y$ . Thus, by definition of  $\Delta$ ,  $\Delta(a) = 0$ , and the statement is proved.

The end of the proof is devoted to proving property (P).

We first study the number of critical points of the projection of  $\text{Zer}(P_\delta, \mathbb{C}^2)$  on the  $X$ -axis close to a critical point of the projection of  $\text{Zer}(P, \mathbb{C}^2)$  to the  $X$ -axis.

- If  $(x, y)$  is a singular point of  $\text{Zer}(P_a, \mathbb{C}^2)$ ,  $(x + \delta y, y)$  is a singular point of  $\text{Zer}(P_{a+\delta}, \mathbb{C}^2)$ .
- If  $(x, y)$  is an ordinary critical point of the projection of  $\text{Zer}(P_a, \mathbb{C}^2)$  on the  $X_1$ -axis,  $(x + \delta y)$  is not a critical point of the projection of  $\text{Zer}(P_{a+\delta}, \mathbb{C}^2)$  on the  $X_1$ -axis. However we are going to prove that there is an ordinary critical point of the projection of  $\text{Zer}(P_{a+\delta}, \mathbb{C}^2)$  on the  $X_1$ -axis of the form  $(x + \delta y + u, y + v)$ ,  $o(u) > 1, o(v) > 0$ .

If  $(x, y)$  is an ordinary critical point of the projection of  $\text{Zer}(P_a, \mathbb{C}^2)$  on the  $X$ -axis, there exists  $r$  such that

$$\begin{aligned} P_a(x, y) = 0, \partial P_a / \partial X(x, y) = b \neq 0, \\ \partial P_a / \partial Y(x, y) = \dots = \partial^{r-1} P_a / \partial Y^{r-1}(x, y) = 0, \partial^r P_a / \partial Y^r(x, y) = r! c \neq 0. \end{aligned}$$

Writing Taylor’s formula for  $P_a$  in the neighborhood of  $(x, y)$ , we have

$$P_a(X, Y) = b(X - x) + c(Y - y)^r + A(X, Y),$$

with  $A(X, Y)$  a linear combination of monomials multiple of  $(X - y)$  or  $(Y - y)^r$ . We consider a new variable  $\varepsilon$  and a point of  $\text{Zer}(P_a, C\langle\varepsilon\rangle^2)$  with coordinates  $x + \xi, y + \varepsilon$ , infinitesimally close to  $x, y$ . In other words we consider a solution  $\xi$  of

$$b\xi + c\varepsilon^r + A(x + \xi, y + \varepsilon) = 0 \tag{11.5}$$

which is infinitesimal. Using the proof of Theorem 2.91, there is a solution of Equation (11.5) of the form  $\xi = -(c/b)\varepsilon^r + w, 0(w) > r$ . Moreover

$$\begin{aligned} \partial P_a / \partial X(x + \xi, y + \varepsilon) &= b + w_1, o(w_1) > 1, \\ \partial P_a / \partial Y(x + \xi, y + \varepsilon) &= cr\varepsilon^{r-1} + w_2, o(w_2) > r - 1. \end{aligned}$$

Thus, if  $(\partial P_a / \partial X(x + \xi', y + \varepsilon'), \partial P_a / \partial Y(x + \xi', y + \varepsilon'))$  is proportional to  $(1, \delta)$ , with  $\varepsilon'$  and  $\xi'$  infinitesimal in  $C\langle\delta\rangle$ , we have

$$\begin{aligned} \varepsilon' &= d\delta^{1/(r-1)} + w_3 \\ \xi' &= -(1/r)d\delta^{r/(r-1)} + w_4 \end{aligned}$$

with  $d^{r-1} = b/cr, o(w_3) > 1/(r - 1), o(w_4) > r/(r - 1)$ . Thus there is a point  $(x + \xi', y + \varepsilon')$  of  $\text{Zer}(P_a, C\langle\delta\rangle^2)$  with gradient proportional to  $(1, \delta)$ , and  $o(\varepsilon') > 0, o(\xi') > 1$ . In other words,  $(x + \delta y + \xi' + \delta\varepsilon', y + \varepsilon')$  is a critical point of the projection of  $\text{Zer}(P_{a+\delta}, C\langle\delta\rangle^2)$  on the  $X$ -axis.

Suppose that  $(x, y_1)$  and  $(x, y_2)$  are two distinct critical point of the projection of  $\text{Zer}(P_a, C^2)$  to the  $X_1$ -axis. According to the preceding discussion, there are two critical point of the projection of  $\text{Zer}(P_{a+\delta}, C^2)$  to the  $X$ -axis, with first coordinates  $x + \delta y_1 + u_1$  and  $x + \delta y_2 + u_2$ , with  $o(u_1) > 1, o(u_2) > 1$ . Note that  $x + \delta y_1 + u_1$  and  $x + \delta y_2 + u_2$  are distinct, since  $y_1 - y_2$  is not infinitesimal. We have proved that the number of distinct roots of  $\text{Disc}(a + \delta, X)$  in  $C\langle\delta\rangle$  is strictly bigger than the number of distinct roots of  $\text{Disc}(a, X)$  in  $C$ .  $\square$

The Topology of a Curve Algorithm can now be described. We perform the computation as if the curve was in generic position. If it is not, the algorithm detects it and then a new linear change of coordinates is performed.

*Algorithm 11.18. [Topology of a Curve]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$  (resp. the field of real numbers  $\mathbb{R}$ ).
- **Input:** a separable polynomial  $P$  in  $D[X, Y]$  of degree  $d$ .
- **Output:** the topology of the curve  $\text{Zer}(P, R^2)$ , described by
  - An integer  $a$  such that  $P_a(X, Y) = P(X + aY, Y)$  is in general position.
  - The number  $r$  of roots of  $\text{Disc}(a, X)$  in  $R$ . We denote these roots by  $x_1 < \dots < x_r$ , and by  $x_0 = -\infty, x_{r+1} = +\infty$ .
  - The number  $m_i$  of roots of  $P_a(x, Y)$  in  $R$  when  $x$  varies on  $(x_i, x_{i+1})$ ,  $i = 0, \dots, r$ .



- The number  $n_i$  of roots of  $P_a(x_i, Y)$  in  $\mathbb{R}$ . We denote these roots by  $y_{i,j}$ ,  $j = 1, \dots, n_i$ .
- An index  $c_i \leq n_i$  such that if  $C_i = (x_i, z_i)$  is the unique critical point of the projection of  $\text{Zer}(P_a, \mathbb{C}^2)$  on the  $X_1$ -axis above  $x_i$ ,  $z_i = y_{i,c_i}$ .

More precisely, the output is  $[a, [m_0, [n_1, c_1], \dots, m_{r-1}, [n_r, c_r], m_r]]$ .

- **Complexity:**  $O(d^{11} \log_2(d))$ , where  $d$  is a bound on the degrees of  $P$ .
- **Procedure:**

- Take  $a := 0$ .
- ( $\star$ ) Define  $P_a(X, Y) := P(X + aX, Y)$ .
- If  $P_a$  is not quasi-monic with respect to  $Y$ , the curve is not in generic position. Take  $a := a + 1$ , go to ( $\star$ ).
- Otherwise, compute the subdiscriminants  $\text{sDisc}(a, X)$  of  $P_a(X, Y)$  (considered as a polynomial in  $Y$ ). Characterize the roots  $x_1, \dots, x_r$  of  $\text{Disc}(a, X)$  using Algorithm 10.14 (Thom encoding) (resp. Algorithm 10.4 (Real Root Isolation)).
- For every  $1 \leq i \leq r$ , determine  $j(i)$  such that

$$\text{Disc}_0(a, x_i) = 0, \dots, \text{sDisc}_{j(i)-1}(a, x_i) = 0, \text{sDisc}_{j(i)}(a, x_i) \neq 0,$$

and compute the sign  $\varepsilon_i$  of  $\text{sDisc}_{j(i)}(a, x_i)$  using Algorithm 10.15 (Sign at the Roots in a real closed field) (resp. Algorithm 10.6 (Sign at a Real Root)).

- Define

$$z_i = -\frac{1}{j(i)} \frac{\text{sDisc}_{j(i), j(i)-1}(a, x_i)}{\text{sDisc}_{j(i)}(a, x_i)}.$$

- Check whether, for every  $i = 1, \dots, r$ ,  $z_i$  is a root of multiplicity  $j(i)$  of  $\text{sDisc}_{j(i)}(x_i, Y)$  using Algorithm 10.15 (Sign at the Roots in a real closed field) (resp. Algorithm 10.6 (Sign at a Real Root)).
- If there exists  $i$  such that  $z_i$  is not a root of multiplicity  $j(i)$  of  $\text{sDisc}_{j(i)}(a, x_i, Y)$ , the curve is not in generic position. Take  $a := a + 1$ , go to ( $\star$ ).
- If for every  $i$ ,  $z_i$  is a root of multiplicity  $j(i)$  of  $\text{sDisc}_{j(i)}(a, x_i, Y)$ , the curve is in generic position.
- For every  $0 \leq i \leq r$  choose an intermediate point  $t_i \in (x_i, x_{i+1})$  (with the convention  $x_0 = -\infty, x_{r+1} = +\infty$ ) using Algorithm 10.18 (Intermediate Points) and Remark 10.77, and compute the number  $m_i$  of roots of  $P_a(t_i, Y)$  evaluating the signs of the signed subresultant coefficients at  $t_i$  using Algorithm 10.11 (Sign Determination) (resp. choose an intermediate rational point  $t_i \in (x_i, x_{i+1})$  using the isolating intervals for the  $x_i$  and compute the number  $m_i$  of roots of  $P_a(t_i, Y)$  using Algorithm 10.4 (Real Root Isolation)).
- Compute, for every  $i \leq r$ ,

$$\begin{aligned} S &= \varepsilon_i(j(i) \text{sDisc}_{j(i)}(a, X) Y + \text{sDisc}_{j(i), j(i)-1}(X) (\partial P_a / \partial Y)) \\ R &= S - \varepsilon_i j(i) d \text{sDisc}_{j(i)}(a, X) P_a \end{aligned}$$

and the list

$$T_i(X) = \text{sRes}(P_a, R).$$

- Evaluate the signs of the elements of  $T_i(X)$  at  $x_i$  to determine the number of real roots of the polynomial  $P_a(x_i, Y)$  which are strictly bigger than  $z_i$  using Proposition 11.22 and Algorithm 10.15 (Sign at the Roots in a real closed field) (resp. Algorithm 10.6 (Sign at a Real Root)).
- Decide whether  $\text{Zer}(P, R^2)$  is empty or a finite number of points.

**Proof of correctness:** By Lemma 4.74, there are only a finite number of values of  $a$  such that  $\deg_{X_2}(Q(a, X, Y)) \neq \deg(Q(a, X, Y))$ . Moreover there are only a finite number of zeros of  $\Delta(X)$ . So by Proposition 11.23, there are only a finite number of values of  $a$  such that the curve is not in generic position. The correctness of the algorithm follows from Proposition 11.21, and the correctness of Algorithm 10.14 (Thom encoding) (resp. Algorithm 10.4 (Real Root Isolation)), Algorithm 10.15 (Sign at the Roots in a real closed field) (resp. Algorithm 10.6 (Sign at a Real Root)), Algorithm 10.18 (Intermediate Points), and Algorithm 10.11 (Sign Determination).  $\square$

**Complexity analysis:** We estimate the complexity of the computation in a general real closed field.

Let  $d$  be the degree of  $P(X, Y)$  a separable polynomial in  $\mathbb{D}[X, Y]$ .

The degree of  $\text{Disc}(U, X)$  is  $O(d^2)$ , and the degree of  $\delta$  is  $O(d^4)$ . So there are at most  $O(d^4)$  values of  $a$  to try.

For each value of  $a$ , we compute the subdiscriminant sequence of  $P_a(X, Y)$  (considered as polynomial in  $Y$ ): this requires  $O(d^2)$  arithmetic operations in  $\mathbb{D}[X]$  by Algorithm 8.21 (Signed subresultant). The degrees of the polynomials in  $X$  produced in the intermediate computations of the algorithm are bounded by  $2d^2$  by Proposition 8.45. The complexity in  $\mathbb{D}$  for this computation is  $O(d^6)$ .

The Thom encoding of the roots of  $\text{Disc}(a, X)$  takes  $O(d^8 \log_2(d))$  arithmetic operations using the complexity analysis of Algorithm 10.14 (Thom encoding). Checking whether  $Q(X, Y)$  is in generic position takes also  $O(d^8 \log_2(d))$  arithmetic operations using the complexity analysis of Algorithm 10.15 (Sign at the Roots in a real closed field), since there are at most  $O(d^2)$  calls to this algorithm with polynomials of degree  $d^2$ .

Since there are at most  $O(d^4)$  values of  $a$  to try, the complexity to reach generic position is  $O(d^{11} \log_2(d))$ .

The choice of the intermediate points, the determination of all  $m_i$ , the determination of  $j(i)$  and  $\varepsilon_i$  takes again  $O(d^8 \log_2(d))$  arithmetic operations using the complexity analysis of Algorithm 10.18 (Intermediate Points) and Algorithm 10.15 (Sign at the Roots in a real closed field).

The computation of one list  $T_i$  by Algorithm 8.21 (Signed subresultant) requires  $O(d^2)$  arithmetic operations in  $D[X]$ . The degree of the polynomials in  $X$  produced in the intermediate computations of the algorithm are bounded by  $O(d^4)$  by Proposition 8.45. The complexity in  $D$  for the computation of the list  $\bar{T}$  is  $O(d^{10})$ .

So the total complexity for computing the various  $T_i$  is  $O(d^{11})$ .

The only remaining computations are the determination of the signs taken by the polynomials of  $T_i$  evaluated at the roots of  $\text{Disc}(a, X)$ . The degrees of the polynomials involved are  $O(d^4)$  and their number is  $O(d^2)$ . This takes  $O(d^{10})$  operations using the complexity analysis of Algorithm 10.15 (Sign at the Roots in a real closed field).

So the total complexity for determining the topology of  $\text{Zer}(P, \mathbb{R}^2)$  is  $O(d^{11} \log_2(d))$ .

If  $D = \mathbb{Z}$  and the coefficients of  $P$  are of bitsizes bounded by  $\tau$ , the coefficients of  $P_a(X, Y)$  are  $\tau + O(d)\nu$  where  $\nu$  is the bitsize of  $d$ , since there are at most  $O(d^4)$  values of  $a$  to try. For each value of  $a$ , the bitsizes of the integers produced in the intermediate computations of the subresultant sequence of  $P_a(X, Y)$  are bounded by  $O(d(\tau + d\nu))$  by the complexity analysis of Algorithm 8.21 (Signed subresultant).

The univariate sign determinations performed in the last steps of the algorithm produce integers of bitsizes  $O(d^2(\tau + d\nu) \log_2(d))$ , according to the complexity analysis of Algorithm 10.14 (Thom encoding), Algorithm 10.15 (Sign at the Roots in a real closed field), Algorithm 10.18 (Intermediate Points). So the maximum bitsizes of the integers in the computation determining the topology of  $\text{Zer}(P, \mathbb{R}^2)$  is  $O(d^2(\tau + d\nu) \log_2(d))$ . □

*Example 11.24.* Let us consider the real algebraic curve defined by the monic polynomial

$$P = 2 Y^3 + (3 X - 3) Y^2 + (3 X^2 - 3 X) Y + X^3.$$

The discriminant of  $P$  (considered as polynomial in  $Y$ ) is

$$-27X^2(X^4 + 6X^2 - 3).$$

whose real roots are given by

$$x_1 = -\sqrt{-3 + 2\sqrt{3}}, \quad x_2 = 0, \quad x_3 = \sqrt{-3 + 2\sqrt{3}}.$$

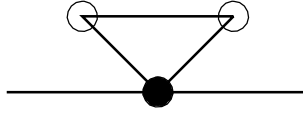
Using the signed subresultant sequence, we determine that above each  $x_{1,i}$ , for  $i = 1, \dots, 3$ ,  $P(x_i, Y)$  has two real roots, and only one of these roots is double. Thus the curve is already in generic position. The multiplicity of the unique critical point  $(x_i, z_i)$  as a root of  $P(x_i, Y)$  is 2, and is given by the equation

$$z_i = -\frac{x_i(x_i^2 + 2x_i - 1)}{2x_i^2 - 2}.$$

We also determine that before  $x_1$  and after  $x_3$ ,  $P(x, Y)$  has only one real root, while between  $x_1$  and  $x_2$ , and between  $x_2$  and  $x_3$   $P(x, Y)$  has three real roots.

For  $i = 1, 3$ , there is one root of  $P(x_i, Y)$  under  $z_i$ , while for  $i = 2$  there is one root of  $P(x_i, Y)$  above  $z_i$ .

Finally, the topology of the curve is given by the following graph



**Fig. 11.3.** Topology of the curve

where the points represent the critical points of the projection on the  $X$ -axis. The white points are critical points of the projection on the  $X$ -axis, while the black one is singular.

The topology can be described by 0 which is the value of  $a$  (the curve was given in generic position) and the list

$$[1, [2, 2], 3, [2, 1], 3, [2, 2], 1]$$

which can be read as follows: there are three roots of  $\text{Disc}(a, X)$ , the number of branches above  $(x_0, x_1)$  is 1, the number of points above  $x_1$  is 2 and the index of the critical point is 2, the number of branches above  $(x_1, x_2)$  is 3, the number of points above  $x_2$  is 2 and the index of the critical point is 1, the number of branches above  $(x_2, x_3)$  is 3, the number of points above  $x_3$  is 2, and the index of the critical point is 2, the number of branches above  $(x_3, x_4)$  is 1.  $\square$

## 11.7 Restricted Elimination

A variant of Algorithm 11.1 (Elimination) will be useful in the last chapters of the book. In this variant, we are interested at the signs of a family of polynomials at the roots of a polynomial, rather than at all the sign conditions of a family of polynomials.

We need some preliminary work.

We first prove a result similar to Theorem 5.12 and Proposition 5.13 which is adapted to the situation we are interested in studying.

**Proposition 11.25.** *Let  $P$  and  $Q$  be in  $\mathbb{R}[X_1, \dots, X_k]$ , and let  $S$  be a semi-algebraically connected semi-algebraic subset of  $\mathbb{R}^{k-1}$  such that  $P$  is not identically 0 on  $S$ , and such that  $\deg(P(x', X_k))$ , the number of distinct roots of  $P$  in  $\mathbb{C}$ , and  $\deg(\gcd(P(x', X_k), Q(x', X_k)))$  are constant over  $S$ . Then there are  $\ell$  continuous semi-algebraic functions  $\xi_1 < \dots < \xi_\ell: S \rightarrow \mathbb{R}$  such that, for every  $x' \in S$ , the set of real roots of  $P(x', X_k)$  is exactly  $\{\xi_1(x'), \dots, \xi_\ell(x')\}$ . Moreover, for  $i = 1, \dots, \ell$ , the multiplicity of the root  $\xi_i(x')$  is constant for  $x' \in S$  and so is the sign of  $Q(\xi_i(x'))$ .*

**Proof:** Let  $a' \in S$  and let  $z_1, \dots, z_j$  be the distinct roots of  $P(a', X_k)$  in  $\mathbb{C}$ . Let  $\mu_i$  (resp.  $\nu_i$ ) be the multiplicity of  $z_i$  as a root of  $P(a', X_k)$  (resp.  $\gcd(P(a', X_k), Q(a', X_k))$ ). The degree of  $P(a', X_k)$  is  $\sum_{i=1}^j \mu_i$  and the degree of  $\gcd(P(a', X_k), Q(a', X_k))$  is  $\sum_{i=1}^j \nu_i$ . Choose  $r > 0$  such that all disks  $D(z_i, r)$  are disjoint.

Using Theorem 5.12 and the fact that the number of distinct complex roots of  $P(x', X_k)$  stays constant over  $S$ , we deduce that there exists a neighborhood  $V$  of  $a'$  in  $S$  such that for every  $x' \in V$ , each disk  $D(z_i, r)$  contains one root of multiplicity  $\mu_i$  of  $P(x', X_k)$ . Since the degree of  $\gcd(P(x', X_k), Q(x', X_k))$  is equal to  $\sum_{i=1}^j \nu_i$ , this gcd must have exactly one root  $\zeta_i$ , of multiplicity  $\nu_i$ , in each disk  $D(z_i, r)$  such that  $\nu_i > 0$ . If  $z_i$  is real,  $\zeta_i$  is real (otherwise, its conjugate  $\bar{\zeta}_i$  would be another root of  $P(x', X_k)$  in  $D(z_i, r)$ ). If  $z_i$  is not real,  $\zeta_i$  is not real, since  $D(z_i, r)$  is disjoint from its image by conjugation. Hence, if  $x' \in V$ ,  $P(x', X_k)$  has the same number of distinct real roots as  $P(a', X_k)$ . Since  $S$  is semi-algebraically connected, the number of distinct real roots of  $P(x', X_k)$  is constant for  $x' \in S$  according to Proposition 3.9. Let  $\ell$  be this number. For  $1 \leq i \leq \ell$ , denote by  $\xi_i: S \rightarrow \mathbb{R}$  the function which sends  $x' \in S$  to the  $i$ -th real root (in increasing order) of  $P(x', X_k)$ . The above argument, with arbitrarily small  $r$ , also shows that the functions  $\xi_i$  are continuous.

It follows from the fact that  $S$  is semi-algebraically connected that each  $\xi_i(x')$  has constant multiplicity as a root of  $P(x', X_k)$  and as a root of  $\gcd(P(x', X_k), Q(x', X_k))$  (cf Proposition 3.9). Moreover, if the multiplicity of  $\xi_i(x')$  as a root of  $\gcd(P(x', X_k), Q(x', X_k))$  is 0, the sign of  $Q$  is fixed on  $\xi_i(x')$ . If  $S$  is described by the formula  $\Theta(X')$ , the graph of  $\xi_i$  is described by the formula

$$\Theta(X') \wedge ((\exists Y_1) \dots (\exists Y_\ell) (Y_1 < \dots < Y_\ell \wedge P(X', Y_1) = 0, \dots, P(X', Y_\ell) = 0)) \wedge ((\forall Y) P(X', Y) = 0 \Rightarrow (Y = Y_1 \vee \dots \vee Y = Y_\ell)) \wedge X_k = Y_i,$$

which shows that  $\xi_i$  is semi-algebraic. □

*Algorithm 11.19. [Restricted Elimination]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a variable  $X_k$ , a polynomial  $P$ , and a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** a finite set  $\text{RElim}_{X_k}(P, \mathcal{P}) \subset D[X_1, \dots, X_{k-1}]$ . The set  $\text{RElim}_{X_k}(P, \mathcal{P})$  is such that the degree of  $P$ , the number of roots of  $P$  in  $\mathbb{R}$ , the number of common roots of  $P$  and  $Q \in \mathcal{P}$  in  $\mathbb{R}$ , and the sign of  $Q \in \mathcal{P}$  at the roots of  $P$  in  $\mathbb{R}$  is fixed on each semi-algebraically connected component of the realization of a sign condition on  $\text{RElim}_{X_k}(P, \mathcal{P})$ .
- **Complexity:**  $s d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $P$  and  $\mathcal{P}$ .

• **Procedure:**

- Place in  $\text{RElim}_{X_k}(P, \mathcal{P})$  the following polynomials when they are not in  $D$ , using Algorithm 8.21 (Signed Subresultant) and Remark 8.50:
  - $\text{sRes}_j(R, \partial R/\partial X_k, R \in \text{Tru}(P), j = 0, \dots, \deg(R) - 2$  (see Definition 1.16).
  - $\text{sRes}_j(\partial R/\partial X_k Q, R)$  for  $Q \in \mathcal{P}, R \in \text{Tru}(P), j = 0, \dots, \deg_{X_k}(R) - 1$ .
  - $\text{lcof}(R)$  for  $R \in \text{Tru}(P)$ .

**Proof of correctness:** The correctness of the Restricted Elimination Algorithm follows from Proposition 11.25. □

**Complexity analysis of Algorithm 11.1:** Consider

$$D[X_1, \dots, X_k] = D[X_1, \dots, X_{k-1}][X_k].$$

There are at most  $d + 1$  polynomials in  $\text{Tru}(P)$  and  $s$  polynomials in  $\mathcal{P}$  so the number of signed subresultant sequences to compute is  $O((s + d)d)$ . Each computation of a signed subresultant sequence costs  $O(d^2)$  arithmetic operations in the integral domain  $D[X_1, \dots, X_{k-1}]$  by the complexity analysis of Algorithm 8.21 (Signed Subresultant). So the complexity is  $O((s + d)d^3)$  in the integral domain  $D[X_1, \dots, X_{k-1}]$ . There are  $O((s + d)d^2)$  polynomials output, of degree bounded by  $2(d^2)$  in  $D[X_1, \dots, X_{k-1}]$ , by Proposition 8.45.

Since each multiplication and exact division of polynomials of degree  $2(d^2)$  in  $k - 1$  variables costs  $O(d)^{4k}$  (see Algorithms 8.5 and 8.6), the complexity in  $D$  is  $s d^{O(k)}$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and output are bounded by  $\tau d^{O(k)}$ , using Proposition 8.46. □

In some phases of our algorithms in the next chapters, we construct points whose coordinates belong to the field of algebraic Puiseux series  $\mathbb{R}\langle\varepsilon\rangle$ . We are going to see that it possible to replace these infinitesimal quantities by sufficiently small elements from the field  $\mathbb{R}$ , using the preceding restricted elimination.

The next proposition makes it possible to replace infinitesimal quantities with sufficiently small elements from the field  $\mathbb{R}$ , using  $\text{RElim}_T$ . For this, we need a bound on the smallest root of a polynomial in terms of its coefficients. Such a bound is given by Proposition 10.3.

We again use the notation introduced in Chapter 10 (Notation 10.1). Given a set of univariate polynomials  $\mathcal{A}$ , we define  $c'(\mathcal{A}) = \min_{Q \in \mathcal{A}} c'(Q)$ .

**Proposition 11.26.** *Let  $f(\varepsilon, T) \in D[\varepsilon, T]$  be a bivariate polynomial,  $\mathcal{L}$  a finite subset of  $D[\varepsilon, T]$ , and  $\sigma$  a sign condition on  $\mathcal{L}$  such that  $f$  has a root  $\bar{t} \in R\langle\varepsilon\rangle$  for which*

$$\bigwedge_{g \in \mathcal{L}} \text{sign}(g(\varepsilon, \bar{t})) = \sigma(g).$$

Then, for any  $v$  in  $\mathbb{R}$ ,  $0 < v < c'(\text{RElim}_T(f, \mathcal{L}))$ , there exists a root  $t$  of  $f(v, T)$  having the same Thom encoding as  $\bar{t}$  and such that

$$\bigwedge_{g \in \mathcal{L}} \text{sign}(g(v, t) = \sigma(g).$$

**Proof:** If  $v < c'(\text{RElim}_T(f, \mathcal{L}))$ , then  $v$  is smaller than the absolute value of all roots of every  $Q$  in  $c'(\text{RElim}_T(f, \mathcal{L}))$  by Proposition 10.3. Hence, by the properties of the output of  $\text{RElim}_T(f, \mathcal{L})$ , the number of roots of the polynomial  $f(v, T)$  as well as the number of its common roots with the polynomials  $g(v, T)$ ,  $g \in \mathcal{L}$ , and the Thom encodings of its roots remain invariant for all  $v$  satisfying  $0 < v < c'(\text{RElim}_T(f, \mathcal{L}))$ .

Since  $\varepsilon < c'(\text{RElim}_T(f, \mathcal{L}))$ , it is clear that for all  $g \in \mathcal{L}$ ,

$$\text{sign}(g(v, t)) = (\text{sign}(g(\varepsilon, \bar{t}))). \quad \square$$

*Algorithm 11.20.* **[Removal of Infinitesimals]**

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $f(\varepsilon, T) \in D[\varepsilon, T]$  and a finite set  $\mathcal{L} \subset D[\varepsilon, T]$ .
- **Output:** a pair  $(a, b)$  of elements of  $D$  such that, for all  $v \in \mathbb{R}$  satisfying  $v \leq a/b$ , the following remains invariant:
  - the number of roots of  $f(v, T)$ , and their Thom encodings,
  - the signs of the polynomials  $g \in \mathcal{L}$  at these roots.
- **Complexity:**  $m d^{O(1)}$ , where  $m$  is a bound on the number of elements of  $\mathcal{L}$  and  $d$  is a bound on the degrees of  $f$  and  $\mathcal{L}$ .
- **Procedure:** Compute  $\text{RElim}_T(f, \mathcal{L}) \subset D[\varepsilon]$  and  $c'(\text{RElim}_T(f, \mathcal{L}))$ . Take  $a$  and  $b$  such that  $a/b = c'(\text{RElim}_T(f, \mathcal{L}))$ .

**Complexity analysis:** According to the complexity analysis of Algorithm 11.19, the complexity is  $m d^{O(1)}$  in  $D$ , since  $k = 2$ .

Note that if  $D = \mathbb{Z}$  and the bitsizes of the coefficients of the polynomials  $f, g$  are bounded by  $\tau$  then  $c'(\text{RElim}_T(f, \mathcal{L}))$  is bounded from below by rational numbers with numerators and denominators of bit-sizes  $\tau d^{O(1)}$ , using the complexity analysis of Algorithm 11.19. In this case, we replace the infinitely small element with a rational number smaller than  $c'(\text{RElim}_T(f, \mathcal{L}))$ . □

*Remark 11.27.* If the coefficients of  $f$  and  $\mathcal{L}$  belong to  $D[w]$  where  $w$  is the root of a polynomial  $h$  of degree at most  $d$  with Thom encoding  $\sigma$ , Algorithm 11.20 (Removal of Infinitesimals) can be easily modified to output a pair  $(a, b)$  of elements of  $D$  such that, for all  $v \in \mathbb{R}$  satisfying  $v \leq a/b$ , the following remains invariant:

- the number of roots of  $f(v, T)$ , and their Thom encodings,
- the signs of the polynomials  $g \in \mathcal{L}$  at these roots.

We just replace  $(a(w), b(w))$  in  $D[w]$  computed by Algorithm 11.20 (Removal of Infinitesimals) by  $(\alpha, \beta)$  where  $\alpha = c'(A)$  and  $\beta = c(B)$  where  $A = \text{Res}_Y(h(Y), T - a(Y))$  (resp.  $B = \text{Res}_Y(h(Y), T - a(Y))$ ). The complexity is clearly  $m d^{O(1)}$ .  $\square$

## 11.8 Bibliographical Notes

The cylindrical decomposition algorithm, due to Collins [45], is the first algorithm for quantifier elimination with a reasonable worst-case time bound. The complexity of the algorithm is polynomial in the degree and number of polynomials. However the complexity is doubly exponential in the number of variables [120]. The former proofs of quantifier elimination [156, 148, 43, 92] were effective, but the complexity of the associated algorithm is not elementary recursive, i.e. is not bounded by a tower of exponents of finite height [120]. The main reason for the improvement in complexity given by the cylindrical decomposition is the use of subresultant coefficients, since, using subresultants, the number of branches in the computation is better controlled.



---

## Polynomial System Solving

This chapter is mainly devoted to algorithms for solving zero-dimensional polynomial systems and give some applications. In Section 12.1, we explain a few results on Gröbner bases. This enables us to decide in Section 12.2 whether a polynomial system is zero-dimensional. We use these results to design various algorithms for zero-dimensional systems, for instance computing the multiplication table for the quotient ring and using the multiplication table to compute information about the solutions of zero-dimensional systems. A special case is treated in details in Section 12.3. In Section 12.4, we define the univariate representations and use trace computations to express the solutions of a zero-dimensional system as rational functions of the roots of a univariate polynomial. In Section 12.5, we explain how to compute the limits of bounded algebraic Puiseux series which are zeros of polynomial systems. In Section 12.6, we introduce the notion of pseudo-critical points and design an algorithm for finding at least one point in every semi-algebraically connected component of a bounded algebraic set, using a variant of the critical point method. In Section 12.8, we describe an algorithm computing the Euler-Poincaré characteristic of an algebraic set.

Throughout this chapter, we assume that  $K$  is an ordered field contained in a real closed field  $R$  and that  $C = R[i]$ .

### 12.1 A Few Results on Gröbner Bases

Throughout Section 12.1 and Section 12.2, we identify a monomial  $X^\alpha \in \mathcal{M}_k$  and the corresponding  $\alpha \in \mathbb{N}^k$ , and we fix a monomial ordering  $<$  on  $\mathcal{M}_k$  (see Definition 4.62). We consider the problem of computing the Gröbner bases of an ideal (see Definition 4.66), using the notation of Section 4.4.1.

It is useful to be able to decide whether a set of polynomials  $\mathcal{G} \subset K[X_1, \dots, X_k]$  is a Gröbner basis. This is done using the notion of an S-polynomial. Given two polynomials  $P_1$  and  $P_2$ , the **S-Polynomial**

of  $P_1$  and  $P_2$  is defined as

$$S(P_1, P_2) = \frac{\text{lt}(P_2)}{g} P_1 - \frac{\text{lt}(P_1)}{g} P_2,$$

where  $g = \text{gcd}(\text{lmon}(P_1), \text{lmon}(P_2))$ . Note that

$$\begin{aligned} \text{lcm}(\text{lmon}(P_1), \text{lmon}(P_2)) &= \frac{\text{lmon}(P_2)}{g} \text{lmon}(P_1) \\ &= \frac{\text{lmon}(P_1)}{g} \text{lmon}(P_2), \\ \text{lmon}(S(P_1, P_2)) &< \text{lcm}(\text{lmon}(P_1), \text{lmon}(P_2)). \end{aligned}$$

**Proposition 12.1. [Termination criterion]** *Let  $\mathcal{G} \subset \mathbb{K}[X_1, \dots, X_k]$  be a finite set such that the leading monomial of any element of  $\mathcal{G}$  is not a multiple of the leading monomial of another element in  $\mathcal{G}$ . Then  $\mathcal{G}$  is a Gröbner basis if and only if the S-polynomial of any pair of polynomials in  $\mathcal{G}$  is reducible to 0 modulo  $\mathcal{G}$ .*

**Proof:** If  $\mathcal{G}$  is a Gröbner basis, for any  $G, H \in \mathcal{G}$ , we have  $S(G, H) \in I(\mathcal{G}, \mathbb{K})$ , and thus  $S(G, H)$  is reducible to 0 modulo  $\mathcal{G}$  by Proposition 4.67.

Conversely, suppose that the S-polynomial of any pair of polynomials in  $\mathcal{G}$  is reducible to 0 modulo  $\mathcal{G}$ . Using Remark 4.65, this implies that for every pair  $G, H \in \mathcal{G}$ ,

$$\exists G_1 \in \mathcal{G} \dots \exists G_s \in \mathcal{G} \quad S(G, H) = \sum_{i=1}^s A_i G_i, \quad (12.1)$$

with  $\text{lmon}(A_i G_i) \leq \text{lmon}(S(G, H))$  for all  $i = 1, \dots, s$ .

Consider  $P \in I(\mathcal{G}, \mathbb{K})$ . We want to prove that  $\text{lmon}(P)$  is a multiple of one of the leading monomials of the elements of  $\mathcal{G}$ . We have  $P = B_1 G_1 + \dots + B_s G_s$  with  $\text{lmon}(P) \leq_{\text{grlex}} \sup \{\text{lmon}(B_i G_i); 1 \leq i \leq s\} = X^\mu$ , and we suppose without loss of generality that  $\text{lmon}(B_1 G_1) = X^\mu$  and  $\text{lcof}(G_i) = 1$ , for all  $i$ .

There are two possibilities:

- $\text{lmon}(P) = X^\mu$ . Then the monomial  $X^\mu$  of  $P$  is a multiple of  $\text{lmon}(G_1)$ , and there is nothing to prove.
- $\text{lmon}(P) <_{\text{grlex}} X^\mu$ . Then the number  $n_\mu$  of  $i$  such that  $\text{lmon}(B_i G_i) = X^\mu$  is at least 2, and we can suppose without loss of generality that  $\text{lmon}(B_2 G_2) = X^\mu$ .

We have

$$\begin{aligned} B_1 G_1 &= b_\beta X^\beta G_1 + F_1 G_1, \\ B_2 G_2 &= c_\gamma X^\gamma G_2 + F_2 G_2, \end{aligned}$$

with  $\text{lmon}(F_1 G_1) < X^\mu$ ,  $\text{lmon}(F_2 G_2) < X^\mu$ .

We necessarily have that  $X^\mu$  is a multiple of

$$X^\sigma = \text{lcm}(\text{lmon}(G_1), \text{lmon}(G_2)).$$

We rewrite  $B_1 G_1 + B_2 G_2$  as follows:

$$\begin{aligned} B_1 G_1 + B_2 G_2 &= (b_\beta + c_\gamma) X^\beta G_1 + c_\gamma (X^\gamma G_2 - X^\beta G_1) + F_1 G_1 + F_2 G_2 \\ &= (b_\beta + c_\gamma) X^\beta G_1 - c_\gamma X^{\mu-\sigma} S(G_1, G_2) + F_1 G_1 + F_2 G_2. \end{aligned}$$

By hypothesis  $S(G_1, G_2) = \sum_{\ell=1}^s A_\ell G_\ell$  with

$$\text{lmon}(A_\ell G_\ell) \leq \text{lmon}(S(G_1, G_2)) <_{\text{grlex}} X^\sigma,$$

and thus

$$\text{lmon}(X^{\mu-\sigma} A_\ell G_\ell) < X^\mu.$$

It follows that we have written

$$P = \bar{B}_1 G_1 + \cdots + \bar{B}_s G_s,$$

and the number of terms  $\bar{B}_i G_i$  with leading monomial  $\mu$  has decreased, or

$$\sup \{\text{lmon}(\bar{B}_i G_i); 1 \leq i \leq s\} < X^\mu.$$

This proves the result by induction on the lexicographical order of the pair  $(\mu, n_\mu)$  since either  $n_\mu$  or  $\nu$  has decreased.  $\square$

The basic idea for computing a Gröbner basis is thus quite simple: add to the original set the reduction of all possible S-polynomials as long as these are not all equal to 0.

*Algorithm 12.1.* **[Buchberger]**

- **Structure:** a field  $K$ .
- **Input:** a finite number of polynomials  $\mathcal{P}$  of  $K[X_1, \dots, X_k]$ .
- **Output:** a Gröbner basis  $\mathcal{G}$  of  $I(\mathcal{P}, K)$ .
- **Procedure:**
  - Reduce  $\mathcal{P}$ , which is done as follows:
    - While there is a polynomial  $P$  in  $\mathcal{P}$  with a monomial  $X^\alpha$  which is a multiple of the leading monomial of an element  $F$  of  $\mathcal{P} \setminus \{P\}$ ,
      - If  $\text{Red}(P, X^\alpha, F) \neq 0$ , update  $\mathcal{P}$  by replacing  $P$  by  $\text{Red}(P, X^\alpha, F)$  in  $\mathcal{P}$ .
      - Otherwise, update  $\mathcal{P}$  by removing  $P$  from  $\mathcal{P}$ .
    - Initialize  $\mathcal{F} := \mathcal{P}$ ,  $\mathcal{C} := \{(P, Q) \mid P \in \mathcal{P}, Q \in \mathcal{P}, P \neq Q\}$ .
    - While  $\mathcal{C} \neq \emptyset$ 
      - Remove a pair  $(P, Q)$  from  $\mathcal{C}$ .
      - Reduce  $S(P, Q)$  modulo  $\mathcal{F}$ , which is done as follows:
        - Initialize  $R := S(P, Q)$

- While there is a monomial of  $R$  that is a multiple of a leading monomial of an element of  $\mathcal{F}$ ,
  - Pick the greatest monomial  $X^\alpha$  of  $R$  which is a multiple of a leading monomial of an element of  $\mathcal{F}$ , and take any  $F \in \mathcal{F}$  such that  $X^\alpha$  is a multiple of  $\text{lmon}(F)$ . Replace  $R$  by  $\text{Red}(R, X^\alpha, F)$ .
- If  $R \neq 0$ , update  $\mathcal{F} := \mathcal{F} \cup \{R\}$ , and reduce  $\mathcal{F}$ :
- While there is a polynomial  $P$  in  $\mathcal{F}$  with a monomial  $X^\alpha$  which is a multiple of the leading monomial of an element  $F$  of  $\mathcal{F} \setminus \{P\}$ ,
  - If  $\text{Red}(P, X^\alpha, F) \neq 0$ , update  $\mathcal{F}$  by replacing  $P$  by  $\text{Red}(P, X^\alpha, F)$  in  $\mathcal{F}$ .
  - Otherwise, update  $\mathcal{F}$  by removing  $P$  from  $\mathcal{F}$ .
- Update  $\mathcal{C} := \mathcal{C} \cup \{(P, Q) \mid P \in \mathcal{F}, Q \in \mathcal{F}, P \neq Q\}$ .
- Output  $\mathcal{G} = \mathcal{F}$ .

**Proof of correctness:** It is first necessary to prove that the algorithm terminates. Denote by  $\mathcal{F}_n$  the value of  $\mathcal{F}$ , obtained after having reduced  $n$  pairs of polynomials. Note that the ideal generated by  $\text{lmon}(\mathcal{F}_{n-1})$  is contained in the ideal generated by  $\text{lmon}(\mathcal{F}_n)$ . Consider the ascending chain of ideals generated by the monomials  $\text{lmon}(\mathcal{F}_n)$ . This ascending chain is stationary by Corollary 4.70 and stops with the ideal generated by  $\text{lmon}(\mathcal{F}_N)$ . This means that for every pair  $(P, Q)$  of elements of  $\mathcal{F}_N$ ,  $S(P, Q)$  is reducible to 0 modulo  $\mathcal{F}_N$ . This ensures that the algorithm terminates. Finally, the leading monomial of any element of  $\mathcal{F}_N$  is not a multiple of the leading monomial of another element in  $\mathcal{F}_N$ , because of the reductions performed at the end of the algorithm. So we conclude that  $\mathcal{G} = \mathcal{F}_N$  is a Gröbner basis of  $I(\mathcal{P}, \mathbf{K})$  by Proposition 12.1.  $\square$

Note that the argument for the termination of Buchberger's algorithm does not provide a bound on the number or degrees of the polynomials output or the number of steps of the algorithm. Thus the complexity analysis of Buchberger's algorithm is a complicated problem which we do not consider here. The reader is referred to [104, 114] for complexity results for the problem of computing Gröbner basis of general polynomial ideals.

*Remark 12.2.* Note that Buchberger's thesis appears in [32] and that a lot of work has been done since then to find out how to compute Gröbner bases efficiently. It is impossible to give in a few words even a vague idea of all the results obtained in this direction (the bibliography [166] contains about 1000 references). In particular, the monomial ordering used plays a key role in the efficiency of the computation, and in many circumstances the reverse lexicographical ordering (Definition 4.63) is a good choice. Modern methods for computing Gröbner bases are often quite different from the original Buchberger's algorithm (see for example [58]). Very efficient Gröbner bases computation can be found at [59].  $\square$

It turns out that certain special polynomial systems are automatically Gröbner bases. This will be very useful for the algorithms to be described in the next chapters.

**Proposition 12.3.** *A polynomial system  $\mathcal{G} = \{X_1^{d_1} + Q_1, \dots, X_\ell^{d_\ell} + Q_\ell\}$  with  $\text{lmon}(Q_i) < X_i^{d_i}$  is a Gröbner basis.*

**Proof:** Let  $P_i = X_i^{d_i} + Q_i$ . We need only prove that  $S(P_i, P_j)$  is reducible to 0 modulo  $\mathcal{G}$ . Clearly  $-Q_i$  is a reduction of  $X_i^{d_i}$  modulo  $\mathcal{G}$ . Note that if  $i \neq j$ , the polynomials  $X_j^{d_j}Q_i$  and  $X_i^{d_i}Q_j$  have no monomials in common. Hence, for  $i \neq j$ ,  $S(P_i, P_j) = X_j^{d_j}Q_i - X_i^{d_i}Q_j$  is reducible to  $X_j^{d_j}Q_i + Q_iQ_j = Q_iP_j$ , which is reducible to 0 modulo  $\mathcal{G}$ . □

Let  $\mathcal{G}$  be a Gröbner basis of an ideal  $I$ . The set of monomials which are multiple of the leading monomials of the polynomials in  $\mathcal{G}$  coincides with the set of leading monomials of elements of  $I$ . The **corners of the staircase** of  $\mathcal{G}$ ,  $\text{Cor}(\mathcal{G})$  are the leading monomials of the polynomials in  $\mathcal{G}$ . They are the minimal elements of the leading monomials of elements of  $I$  for the partial order of divisibility among monomials. The **orthant** generated by a corner  $X^\alpha$  is  $\alpha + \mathbb{N}^k$  and consists of multiples of  $X^\alpha$ . The set of leading monomials of elements of  $I$  is the union of the orthants generated by the corners. The set of **monomials under the staircase** for  $\mathcal{G}$ ,  $\text{Mon}(\mathcal{G})$ , is the set of monomials that do not belong to the set of leading monomials of elements of  $I$ . The **border of the staircase** of  $\mathcal{G}$  is the set  $\text{Bor}(\mathcal{G})$  of monomials which are leading monomials of elements of  $I$  and which are obtained by multiplying a monomial of  $\text{Mon}(\mathcal{G})$  by a variable.

*Example 12.4.* We illustrate these notions using an example in the plane. Let  $k = 2$  and the monomial ordering be the reverse lexicographical ordering. Consider  $\mathcal{G} := \{X_1^2 + 2X_2^2, X_2^4, X_1X_2^2\}$ . It is easy to check that  $\mathcal{G}$  is a Gröbner basis. The set of leading monomials of elements of  $\text{Ideal}(\mathcal{G}, \mathbb{K})$  is the set of multiples of  $X_2^4$ ,  $X_1X_2^2$  and  $X_1^2$ . It is the union of the three orthants generated by the corners  $X_2^4$ ,  $X_1X_2^2$ ,  $X_1^2$ . The set of monomials under the staircase is the finite set of monomials

$$\text{Mon}(\mathcal{G}) = \{1, X_2, X_2^2, X_2^3, X_1, X_1X_2\}.$$

The border of  $\mathcal{G}$  is

$$\text{Bor}(\mathcal{G}) = \{X_1^2, X_1^2X_2, X_1X_2^2, X_1X_2^3, X_2^4\}.$$

Here is the corresponding picture (the big black points are the corners, the other black points are the other elements of the border, the white points are under the staircase).

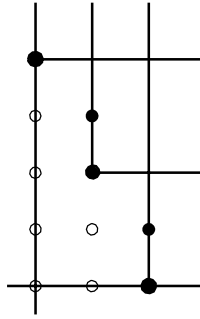


Fig. 12.1. Staircase

□

**Proposition 12.5.** *Let  $A = \mathbb{K}[X_1, \dots, X_k]/I$  for an ideal  $I$  and let  $\mathcal{G}$  be a Gröbner basis of  $I$ . The monomials in  $\text{Mon}(\mathcal{G})$  constitute a basis of the  $\mathbb{K}$ -vector space  $A$ .*

**Proof:** It is clear that any element of  $\mathbb{K}[X_1, \dots, X_k]$  is reducible modulo  $\mathcal{G}$  to a linear combination of monomials in  $\text{Mon}(\mathcal{G})$ . Conversely, a non-zero linear combination of monomials in  $\text{Mon}(\mathcal{G})$  cannot be reduced to 0 by  $\mathcal{G}$ , thus does not belong to  $I$ , and hence is not zero in  $A$ . □

Let  $\mathcal{G}$  be a Gröbner basis. The **normal form** of  $P \in \mathbb{K}[X_1, \dots, X_k]$  modulo  $\mathcal{G}$ , denoted  $\text{NF}(P)$ , is a linear combination  $Q$  of monomials in  $\text{Mon}(\mathcal{G})$  such that  $P = Q \pmod{I(\mathcal{G}, \mathbb{K})}$ . Such a linear combination is unique by Proposition 12.5. Note that  $\text{NF}$  is a linear mapping from the  $\mathbb{K}$ -vector space  $\mathbb{K}[X_1, \dots, X_k]$  to the  $\mathbb{K}$ -vector space  $\mathbb{K}[X_1, \dots, X_k]/I(\mathcal{G}, \mathbb{K})$ .

*Algorithm 12.2.* [Normal Form]

- **Structure:** a field  $\mathbb{K}$ .
- **Input:** a Gröbner basis  $\mathcal{G}$  and a polynomial  $P \in \mathbb{K}[X_1, \dots, X_k]$ .
- **Output:** the normal form  $\text{NF}(P)$  of  $P$  modulo  $\mathcal{G}$ .
- **Procedure:**
  - Initialize  $Q := P$ .
  - While there are monomials of  $Q$  that are not in  $\text{Mon}(\mathcal{G})$ ,
    - Pick the greatest monomial  $X^\alpha$  of  $Q$  which is not in  $\text{Mon}(\mathcal{G})$ . Take any  $G \in \mathcal{G}$  such that  $X^\alpha$  is a multiple of  $\text{lmon}(G)$ . Update  $Q$ , replacing it by  $\text{Red}(Q, X^\alpha, G)$ .
  - Output  $\text{NF}(P) := Q$ .

*Example 12.6.* Returning to example 12.4 the normal form of a polynomial in  $X_1$  and  $X_2$ , is a linear combination of elements in  $\text{Mon}(\mathcal{G})$ . For example, the normal form of  $X_1^3 + X_1 + X_2$  is obtained by computing

$$\begin{aligned} \text{Red}(X_1^3 + X_1 + X_2, X_1^3, X_1^2 + 2X_2^2) &= -2X_1X_2^2 + X_1 + X_2, \\ \text{Red}(-2X_1X_2^2 + X_1 + X_2, X_1X_2^2, X_1X_2^2) &= X_1 + X_2. \end{aligned}$$

□

## 12.2 Multiplication Tables

An important property of a Gröbner basis is that it can be used to characterize zero-dimensional polynomial systems.

**Proposition 12.7.** *The finite set  $\mathcal{P} \subset K[X_1, \dots, X_k]$  is a zero-dimensional polynomial system if and only if any Gröbner basis of  $\text{Ideal}(\mathcal{P}, K)$  contains a polynomial with leading monomial  $X_i^{d_i}$  for each  $i, 1 \leq i \leq k$ .*

**Proof:** It is clear that the number of monomials under the staircase is finite if and only any Gröbner basis contains, for each  $i, 1 \leq i \leq k$ , a polynomial with leading monomial  $X_i^{d_i}$ . We conclude by applying Proposition 12.5 and Theorem 4.85. □

As a consequence:

**Corollary 12.8.** *If  $\mathcal{P} \subset K[X_1, \dots, X_k]$  is a zero-dimensional polynomial system and  $\mathcal{G}$  is its Gröbner basis,  $\#(\text{Zer}(\mathcal{P}, C^k)) \leq \#(\text{Mon}(\mathcal{G}))$ .*

**Proof:** Use Proposition 12.5 and Theorem 4.85. □

**Corollary 12.9.** *A polynomial system  $\mathcal{P} = \{X_1^{d_1} + Q_1, \dots, X_k^{d_k} + Q_k\}$  with  $\text{lmon}(Q_i) <_{\text{grlex}} X_i^{d_i}$  has a finite number of solutions.*

**Proof:** Apply Proposition 12.3 and Proposition 12.7. □

Suppose that  $\mathcal{P} \subset K[X_1, \dots, X_k]$  is a zero-dimensional polynomial system and  $\mathcal{B}$  is a basis of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ . The **multiplication table** of  $A$  in  $\mathcal{B}$  is, for every pair of elements  $a$  and  $b$  in  $\mathcal{B}$ , the expression of their product in  $A$  as a linear combination of elements in  $\mathcal{B}$ . Note that it happens often that  $a = a' b'$ , with  $a, b, a', b'$  in  $\mathcal{B}$ . So in order to avoid repetitions we define

$$\text{Tab}(\mathcal{B}) = \{a b \mid a, b \text{ in } \mathcal{B}\}.$$

The **size of the multiplication table**, denoted by  $T$ , is the number of elements of  $\text{Tab}(\mathcal{B})$ . The number  $T$  can be significantly smaller than  $N^2$  where  $N$  is the number of monomials in  $\text{Tab}(\mathcal{B})$ . For example when  $k=1$ , there are at most  $2N$  monomials in the multiplication table, taking  $\mathcal{B} = \{1, X, \dots, X^{N-1}\}$ . The coefficients  $\lambda_{c,d}$  such that, in  $A$ ,

$$c = \sum_{d \in \mathcal{B}} \lambda_{c,d} d,$$

with  $c \in \text{Tab}(\mathcal{B})$ , are the **entries of the multiplication table**  $\text{Mat}(\mathcal{B})$ .

We explain now how to compute the multiplication table in the basis  $\text{Mon}(\mathcal{G})$ .

*Algorithm 12.3. [Multiplication Table]*

- **Structure:** a field  $K$ .
- **Input:** a zero-dimensional polynomial system  $\mathcal{P} \subset K[X_1, \dots, X_k]$ , together with a Gröbner basis  $\mathcal{G}$  of  $I = \text{Ideal}(\mathcal{P}, K)$ .
- **Output:** the multiplication table  $\text{Mat}(\text{Mon}(\mathcal{G}))$  of  $A = K[X_1, \dots, X_k]/I$  in the basis  $\text{Mon}(\mathcal{G})$ . More precisely, for every  $c \in \text{Tab}(\text{Mon}(\mathcal{G}))$ , the coefficients  $\lambda_{c,a}$ ,  $a \in \text{Mon}(\mathcal{G})$ , such that  $c = \sum_{a \in \text{Mon}(\mathcal{G})} \lambda_{c,a} a$ .
- **Complexity:**  $O(k N^3 + T N^2)$ , where  $N$  is the number of elements of  $\text{Mon}(\mathcal{G})$  and  $T$  is the number of elements of  $\text{Tab}(\text{Mon}(\mathcal{G}))$ .
- **Procedure:**

– Step 1: For every  $X^\alpha \in \text{Bor}(\mathcal{G})$  in increasing order according to  $<$  compute  $\text{NF}(X^\alpha) = \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \mu_{\alpha,\delta} X^\delta$  as follows:

– If  $X^\alpha \in \text{Cor}(\mathcal{G})$ , and  $G \in \mathcal{G}$  is such that  $\text{lmon}(G) = X^\alpha$ ,

$$\text{NF}(X^\alpha) := G - X^\alpha.$$

– If  $X^\alpha \notin \text{Cor}(\mathcal{G})$ , and  $X^\alpha = X_j X^\beta$  for some  $j = 1, \dots, k$ , with  $X^\beta \in \text{Bor}(\mathcal{G})$ , define

$$\text{NF}(X^\alpha) := \sum_{(X^\gamma, X^\delta) \in \text{Mon}(\mathcal{G})^2} \mu_{\beta,\gamma} \mu_{\gamma',\delta} X^\delta,$$

with  $X_j X^\gamma = X^{\gamma'}$ , and

$$\text{NF}(X^\beta) = \sum_{X^\gamma \in \text{Mon}(\mathcal{G})} \mu_{\beta,\gamma} X^\gamma,$$

$$\text{NF}(X^{\gamma'}) = \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \mu_{\gamma',\delta} X^\delta.$$

– Step 2: Construct the matrices  $M_1, \dots, M_k$  corresponding to multiplication by  $X_1, \dots, X_k$ , expressed in the basis  $\text{Mon}(\mathcal{G})$ , using the normal forms of elements of  $\text{Bor}(\mathcal{G})$  already computed.

– Step 3: For every  $X^\alpha \in \text{Tab}(\mathcal{G}) \setminus (\text{Mon}(\mathcal{G}) \cup \text{Bor}(\mathcal{G}))$  in increasing order according to  $<$  compute  $\text{NF}(X^\alpha)$  as follows: since  $X^\alpha = X_j X^\beta$ , for some  $j = 1, \dots, k$ , compute the vector  $\text{NF}(X^\alpha) = M_j \cdot \text{NF}(X^\beta)$ .



**Proof of correctness:** Note that

$$\begin{aligned} \text{NF}(X^\alpha) &= \text{NF}(X_j \text{NF}(X^\beta)) \\ &= \sum_{X^\gamma \in \text{Mon}(\mathcal{G})} \mu_{\beta, \gamma} X_j X^\gamma \\ &= \sum_{(X^\gamma, X^\delta) \in \text{Mon}(\mathcal{G})^2} \mu_{\beta, \gamma} \mu_{\gamma', \delta} X^\delta, \end{aligned}$$

The only thing which remains to prove is that

$$\begin{aligned} \text{NF}(X^\beta) &= \sum_{X^\gamma \in \text{Mon}(\mathcal{G})} \mu_{\beta, \gamma} X^\gamma, \\ \text{NF}(X^{\gamma'}) &= \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \mu_{\gamma', \delta} X^\delta, \end{aligned}$$

have been computed before, which follows from  $X^\beta < X^\alpha$  and  $X^{\gamma'} < X^\alpha$ .  $\square$

**Complexity analysis:**

There are at most  $kN$  elements in  $\text{Bor}(\mathcal{G})$ . For each element of  $\text{Bor}(\mathcal{G})$ , the computation of its normal form takes  $O(N^2)$  operations. Thus, the complexity of Step 1 is  $O(kN^3)$ .

In Step 3 there is a matrix vector multiplication to perform for each element on  $\text{Tab}(\mathcal{G})$ , so the total cost of Step 3 is  $O(TN^2)$ .  $\square$

Once a multiplication table is known, it is easy to perform arithmetic operations in the quotient ring  $A = K[X_1, \dots, X_k]/I$ , and to estimate the complexity of these computations. All the arithmetic operations to perform take place in the ring generated by the entries of the multiplication table and the coefficients of the elements we want to add or multiply.

*Algorithm 12.4.* **[Zero-dimensional Arithmetic Operations]**

- **Structure:** a ring  $D$  contained in a field  $K$ .
- **Input:**
  - a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ ,
  - a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ , such that the multiplication table  $\text{Mat}(\mathcal{B})$  has entries in  $D$ ,
  - two elements  $f, g \in A$ , given as a linear combination of elements of  $\mathcal{B}$  with coefficients in  $D$ .
- **Output:**  $f + g$  and  $fg$  in  $A$ , specified by linear combinations with coefficients in  $D$  of elements of  $\mathcal{B}$ .
- **Complexity:**  $O(N)$  for the addition,  $O(TN^2)$  for the multiplication, where  $N$  is the number of elements of  $\mathcal{B}$  and  $T$  is the number of elements of  $\text{Tab}(\mathcal{B})$ .
- **Procedure:**
  - Let  $f = \sum_{a \in \mathcal{B}} f_a a$ ,  $g = \sum_{a \in \mathcal{B}} g_a a$ .
  - Define  $f + g := \sum_{\alpha \in \mathcal{B}} (a_\alpha + b_\alpha) \alpha$ .

- If  $c = \sum_{d \in \mathcal{B}} \lambda_{c,d} d$ , for  $c \in \text{Tab}(\mathcal{B})$ , define

$$fg := \sum_{d \in \mathcal{B}} \sum_{c \in \text{Tab}(\mathcal{B})} \sum_{\substack{a \in \mathcal{B}, b \in \mathcal{B} \\ ab=c}} (f_a g_b) \lambda_{c,d} d.$$

**Complexity analysis:** It is clear that the complexity of the addition is  $O(N)$ , while the complexity of the multiplication is  $O(TN^2)$ .  $\square$

Lots of information can be obtained about the solutions of a zero-dimensional system once a multiplication table is known. We can compute traces, number of distinct zeroes and perform sign determination.

We use the notation introduced in Chapter 4 (Notation 4.95). For  $f \in A$  we denote by  $L_f: A \rightarrow A$  the linear map defined by  $L_f(g) = fg$  for  $g \in A$ . We can compute the trace of  $L_f$ .

*Algorithm 12.5. [Trace]*

- **Structure:** a ring  $D$  contained in a field  $K$ .
- **Input:**
  - a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ ,
  - a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ , such that the multiplication table  $\text{Mat}(\mathcal{B})$  has entries in  $D$ ,
  - an element  $f \in A$ , given as a linear combination of elements of  $\mathcal{B}$  with coefficients in  $D$ .
- **Output:** the trace of the linear map  $L_f$ .
- **Complexity:**  $O(N^2)$ , where  $N$  is the number of elements of  $\mathcal{B}$ .
- **Procedure:** Let  $f = \sum_{a \in \mathcal{B}} f_a a$ . Compute

$$\text{Tr}(L_f) := \sum_{b \in \mathcal{B}} \sum_{a \in \mathcal{B}} f_a \lambda_{ab,b}.$$

**Proof of correctness:** Since  $fb = \sum_{c \in \mathcal{B}} \sum_{a \in \mathcal{B}} f_a \lambda_{ab,c} c$ , the entry of the matrix of  $L_f$  corresponding to the elements  $b, c$  of the basis is  $\sum_{a \in \mathcal{B}} f_a \lambda_{ab,c}$ . The trace of  $L_f$  is obtained by summing the diagonal terms.  $\square$

**Complexity analysis:** The number of arithmetic operations in  $D$  is clearly  $O(N^2)$ .  $\square$

We can also compute the number of distinct zeroes.

*Algorithm 12.6. [Number of Distinct Zeros]*

- **Structure:** an integral domain  $D$  contained in a field  $K$ .
- **Input:**
  - a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ ,
  - a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ , such that the multiplication table  $\text{Mat}(\mathcal{B})$  has entries in  $D$ ,

- **Output:**  $n = \#\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .
- **Complexity:**  $O(TN^2)$ , where  $N$  is the number of elements of  $\mathcal{B}$  and  $T$  is the number of elements of  $\text{Tab}(\mathcal{B})$ .
- **Procedure:** Apply Algorithm 12.5 (Trace) to the maps  $L_c$  for every element  $c$  of  $\text{Tab}(\mathcal{B})$ . Then compute the rank of the matrix with entries  $[\text{Tr}(L_{ab})]$  for every  $a, b$  in  $\mathcal{B}$ , using Remark 8.18.

**Proof of correctness:** The matrix with entries  $[\text{Tr}(L_{ab})]$  is the matrix of  $\text{Her}(\mathcal{P})$  in the basis  $\mathcal{B}$ . Hence its rank is equal to  $\#\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  by part a) of Theorem 4.100 (Multivariate Hermite). □

**Complexity analysis:** Let  $N$  be the dimension of the  $\mathbb{K}$ -vector space  $A$ . There are  $T$  trace computations to perform, and then a rank computation for a matrix of size  $N$ . The number of arithmetic operations in  $\mathbb{D}$  is thus  $O(TN^2)$  using the complexity analysis of Algorithm 12.5 (Trace), Algorithm 8.16 (Jordan-Bareiss’s method) and Remark 8.18. □

We can also compute Tarski-queries.

*Algorithm 12.7.* **[Multivariate Tarski-query]**

- **Structure:** an integral domain  $\mathbb{D}$  contained in an ordered field  $\mathbb{K}$ .
- **Input:**
  - a zero-dimensional polynomial system  $\mathcal{P} \subset \mathbb{D}[X_1, \dots, X_k]$ ,
  - a basis  $\mathcal{B}$  of  $A = \mathbb{K}[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, \mathbb{K})$ , such that the multiplication table  $\text{Mat}(\mathcal{B})$  has entries in  $\mathbb{D}$ ,
  - an element  $Q \in A$ , given as a linear combination of elements of  $\mathcal{B}$  with coefficients in  $\mathbb{D}$ .

given as a linear combination of elements of  $\mathcal{B}$  with coefficients in  $\mathbb{D}$ .
- **Output:** the Tarski-query  $\text{TaQ}(Q, Z)$  with  $Z = \text{Zer}(\mathcal{P}, \mathbb{R}^k)$ .
- **Complexity:**  $O(TN^2)$ , where  $N$  is the number of elements of  $\mathcal{B}$  and  $T$  is the number of elements of  $\text{Tab}(\mathcal{B})$ .
- **Procedure:** Apply the Trace Algorithm to the maps  $L_{Qc}$  for every  $c \in \text{Tab}(\mathcal{B})$ . Then compute the signature of the Hermite quadratic form associated to  $Q$  by Algorithm 8.18 (Signature through Descartes).

**Proof of correctness:** The correctness of the Multivariate Tarski-query Algorithm follows from Theorem 4.100 (Multivariate Hermite). □

**Complexity analysis:** There are  $T$  traces to compute, and a signature computation to perform. The number of arithmetic operations in  $\mathbb{D}$  is  $O(TN^2)$ , given the complexity analyses of Algorithms 12.5 (Trace) and 8.18 (Signature through Descartes). □

Algorithm 10.11 (Sign Determination) can be performed with Algorithm 12.7 (Multivariate Tarski-query) as a blackbox.

*Algorithm 12.8. [Multivariate Sign Determination]*

- **Structure:** an integral domain  $D$  contained in a field  $K$ .
- **Input:**
  - a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ ,
  - a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ , such that the multiplication table  $\text{Mat}(\mathcal{B})$  has entries in  $D$ ,
  - a list  $\mathcal{Q} \subset A$ , given as a linear combinations of elements of  $\mathcal{B}$  with coefficients in  $D$ .

Denote by  $Z = \text{Zer}(\mathcal{P}, \mathbb{R}^k)$ .

- **Output:** the list of sign conditions realized by  $\mathcal{Q}$  on  $Z$ .
- **Complexity:**  $s T N^3$ , where  $N$  is the number of elements of  $\mathcal{B}$ ,  $T$  is the number of elements of  $\text{Tab}(\mathcal{B})$  and  $s$  the number of elements of  $\mathcal{Q}$ .
- **Procedure:** Perform Algorithm 10.11 (Sign Determination) for  $Z$ , using Algorithm 12.7 (Multivariate Tarski-query) as a Tarski-query black box, the products of polynomials in  $\mathcal{Q}$  being computed in  $A$  using Algorithm 12.4 (Zero-dimensional Arithmetic Operations).

**Complexity analysis:** Let  $N$  be the dimension of the  $K$ -vector space  $A$  and  $s$  the cardinality of  $\mathcal{Q}$ . Note that  $\#(Z) \leq N$  by Theorem 4.85. According to the complexity of Algorithm 10.11 (Sign Determination), the number of calls to the Tarski-query black box is bounded by  $1 + 2sN$ . For each call to Algorithm 12.7 (Multivariate Tarski-query), a multiplication has to be performed by Algorithm 12.4 (Zero-dimensional Arithmetic Operations). The complexity is thus  $O(s T N^3)$ , using the complexity of Algorithm 12.7 (Multivariate Tarski-query).  $\square$

## 12.3 Special Multiplication Table

We now study a very special case of zero-dimensional system, where we start from a Gröbner basis with a very specific structure, used in the later chapters. In this section the only monomial ordering we consider is the graded lexicographical ordering (see Definition 2.15).

Let  $D$  be an integral domain. A Gröbner basis  $\mathcal{G}$  is **special** if it is of the form

$$\mathcal{G} = \{b_1 X_1^{d_1} + Q_1, \dots, b_k X_k^{d_k} + Q_k\} \subset D[X_1, \dots, X_k]$$

with  $\deg(Q_i) < d_i$ ,  $\deg_{X_j}(Q_i) < d_j$ ,  $j \neq i$ ,  $d_1 \geq \dots \geq d_k$ ,  $b_j \neq 0$ . According to Proposition 12.3, a special Gröbner basis  $\mathcal{G}$  is a Gröbner basis of  $\text{Ideal}(\mathcal{G}, K)$  for the graded lexicographical ordering on monomials. The monomials under the staircase  $\text{Mon}(\mathcal{G})$  are the monomials  $X^\alpha = X_1^{\alpha_1} \dots X_k^{\alpha_k}$  with  $\alpha_i < d_i$ . Note that, for every  $i$ ,  $Q_i$  is a linear combination of  $\text{Mon}(\mathcal{G})$ . Thus,  $\text{NF}(b_i X_i^{d_i}) = -Q_i$ . The border of the staircase  $\text{Bor}(\mathcal{G})$  is the set of monomials such that  $\alpha_i = d_i$  for some  $i \in \{1, \dots, k\}$  and  $\alpha_j < d_j$  for all  $j \neq i$ .

We adapt Algorithm 12.3 (Multiplication Table) to this special case, the main change being that we manage to obtain coefficients in  $D$ . In order to be able to make the reduction inside the ring  $D$ , it is useful to multiply monomials in advance by a convenient product of the leading coefficients of the polynomials in  $\mathcal{G}$ . This is the reason why we introduce the following notation.

**Notation 12.10. [Multiplication table]** For  $\alpha \in \mathbb{N}^k$ , let  $|\alpha| = \alpha_1 + \dots + \alpha_k$  denote the total degree of  $X^\alpha$ . Let  $b$  be a common multiple of  $b_1, \dots, b_k$  in  $D$ , i.e. for every  $j = 1, \dots, k$ ,  $b = b_j \bar{b}_j$ ,  $\bar{b}_j \in D$ . Denote by  $\overline{\text{Mon}}(\mathcal{G})$  the set of  $b^{|\alpha|} X^\alpha$  such that  $X^\alpha \in \text{Mon}(\mathcal{G})$  and by  $\overline{\text{Bor}}(\mathcal{G})$  the set of  $b^{|\alpha|} X^\alpha$  such that  $X^\alpha \in \text{Bor}(\mathcal{G})$ .  $\square$

We adapt Algorithm 12.3 (Multiplication Table) to this special case, the main change being that we manage to obtain coefficients in  $D$ .

*Algorithm 12.9. [Special Multiplication Table]*

- **Structure:** a ring  $D$  contained in a field  $K$ .
- **Input:** a special Gröbner basis

$$\mathcal{G} = \{b_1 X_1^{d_1} + Q_1, \dots, b_k X_k^{d_k} + Q_k\} \subset D[X_1, \dots, X_k],$$

$b$  a common multiple of  $b_1, \dots, b_k$  in  $D$ , with, for every  $j$  from 1 to  $k$ ,  $b = b_j \bar{b}_j$ .

- **Output:** the multiplication table of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$  in the basis  $\overline{\text{Mon}}(\mathcal{G})$ . The multiplication table consists of polynomials in  $D[X_1, \dots, X_k]$ .
- **Complexity:**  $O(2^k (d_1 \dots d_k)^3)$ .
- **Procedure:**

– Step 1: For every  $b^{|\alpha|} X^\alpha \in \overline{\text{Bor}}(\mathcal{G})$  in increasing order according to  $<_{\text{deglex}}$  compute  $\text{NF}(b^{|\alpha|} X^\alpha) = \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \lambda_{\alpha, \delta} b^{|\alpha|} X^\delta$  as follows:

- If  $X^\alpha = X_i^{d_i}$ ,  $\text{NF}(b^{|\alpha|} X^\alpha) := -b^{d_i-1} \bar{b}_i Q_i$ .
- Else, if  $X^\alpha = X_j X^\beta$  for some  $j = 1, \dots, k$ , with  $X^\beta \in \text{Bor}(\mathcal{G})$ , define

$$\text{NF}(b^{|\alpha|} X^\alpha) = \sum_{(X^\gamma, X^\delta) \in \text{Mon}(\mathcal{G})^2} \lambda_{\beta, \gamma} \lambda_{\gamma', \delta} b^{|\gamma|} X^\delta,$$

with  $X_j X^\gamma = X^{\gamma'}$ , and

$$\begin{aligned} \text{NF}(b^{|\beta|} X^\beta) &= \sum_{X^\gamma \in \text{Mon}(\mathcal{G})} \mu_{\beta, \gamma} b^{|\gamma|} X^\gamma, \\ \text{NF}(b^{|\gamma'|} X^{\gamma'}) &= \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \mu_{\gamma', \delta} b^{|\delta|} X^\delta. \end{aligned}$$

- Step 2: Construct the matrices  $M'_1, \dots, M'_k$  corresponding to multiplication by  $b X_1, \dots, X b_k$ , expressed in the basis  $\overline{\text{Mon}}(\mathcal{G})$ , using the normal forms of elements of  $\overline{\text{Bor}}(\mathcal{G})$  already computed.
- Step 3: For every  $X^\alpha \in \text{Tab}(\mathcal{G}) \setminus (\text{Mon}(\mathcal{G}) \cup \text{Bor}(\mathcal{G}))$  in increasing order according to  $<_{\text{deglex}}$  compute  $\text{NF}(b^{|\alpha|} X^\alpha)$  as follows: since  $X^\alpha = X_j X^\beta$ , for some  $j = 1, \dots, k$ , compute the vector

$$\text{NF}(b^{|\alpha|} X^\alpha) = M'_j \cdot \text{NF}(b^{|\beta|} X^\beta).$$

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 12.3 (Multiplication table). The fact that the multiplication table consists of polynomials with coefficients in  $D$  is clear.  $\square$

**Complexity analysis:** Follows from the complexity analysis of Algorithm 12.3 (Multiplication Table), noting that the number of elements of  $\text{Mon}(\mathcal{G})$  is  $d_1 \cdots d_k$ , and that the number of elements in  $\text{Tab}(\mathcal{G})$  is bounded by  $2^k d_1 \cdots d_k$ .  $\square$

It will be necessary in several subsequent algorithms to perform the same computation with parameters. The following paragraphs are quite technical but it seems to be unfortunately unavoidable. Let  $Y = Y_1, \dots, Y_\ell$ . We say that  $\mathcal{G}(Y)$  is a **parametrized special Gröbner basis** if it is of the form

$$\mathcal{G}(Y) = \{b_1(Y) X_1^{d_1} + Q_1(Y, X), \dots, b_k(Y) X_k^{d_k} + Q_k(Y, X)\}$$

with  $\deg(Q_i) < d_i$ , where  $\deg$  is the total degree with respect the variables  $X_1, \dots, X_k$ ,  $\deg_{X_j}(Q_i) < d_i$ ,  $j \neq i$ ,  $d_1 \geq \dots \geq d_k$ ,  $b_j \neq 0 \in D[Y]$  for  $0 \leq j \leq k$ . Let  $b(Y)$  a common multiple of  $b_1(Y), \dots, b_k(Y)$  in  $D[Y]$ , with  $b(Y) = b_j(Y) \bar{b}_j(Y)$ . Then, for any  $y \in C^\ell$  such that  $b(y) \neq 0$ ,

$$\mathcal{G}(y) = \{b_1(y) X_1^{d_1} + Q_1(y, X), \dots, b_k(y) X_k^{d_k} + Q_k(y, X)\} \subset C[X_1, \dots, X_k].$$

is a special Gröbner basis. Define  $\overline{\text{Mon}}(\mathcal{G})$  as the set of elements  $b(Y)^{|\alpha|} X^\alpha = b(Y)^{|\alpha|} X_1^{\alpha_1} \cdots X_k^{\alpha_k}$  with  $\alpha_i < d_i$  and  $\overline{\text{Bor}}(\mathcal{G})$  as the set of elements  $b(Y)^{|\alpha|} X^\alpha$  such that  $\alpha_i = d_i$  for some  $i \in \{1, \dots, k\}$  and  $\alpha_i \leq d_i$  for any  $i \in \{1, \dots, k\}$ .

*Algorithm 12.10.* **[Parametrized Special Multiplication Table]**

- **Structure:** a ring  $D$  contained in a field  $K$ .
- **Input:** a parametrized special Gröbner basis

$$\mathcal{G} = \{b_1(Y) X_1^{d_1} + Q_1(Y, X), \dots, b_k(Y) X_k^{d_k} + Q_k(Y, X)\} \subset D[Y][X_1, \dots, X_k]$$

with  $Y = (Y_1, \dots, Y_\ell)$ , and  $b(Y)$  a common multiple of  $b_1(Y), \dots, b_k(Y)$  in  $D[Y]$ , with  $b(Y) = b_j(Y) \bar{b}_j(Y)$ .

- **Output:** a parametrized multiplication table in the basis  $\overline{\text{Mon}}(\mathcal{G})$ : i.e. for any two monomials  $b(Y)^{|\alpha|} X^\alpha$  and  $b(Y)^{|\beta|} X^\beta$  in  $\overline{\text{Mon}}(\mathcal{G})$  a linear combination  $\text{NF}(b(Y)^{|\alpha|+|\beta|} X^\alpha X^\beta)(Y)$  of monomials in  $\overline{\text{Mon}}(\mathcal{G})$  with coefficients in  $D[Y]$  such that for every  $y \in C^\ell$  such that  $b(y) \neq 0$ , the polynomial  $\text{NF}(b(Y)^{|\alpha|+|\beta|} X^\alpha X^\beta)(y)$  is the normal form of  $b(y)^{|\alpha|} X^\alpha b(y)^{|\beta|} X^\beta$  modulo  $\mathcal{G}(y)$ .
- **Complexity:**  $(d_1 \cdots d_k \lambda)^{O(\ell)}$ .
- **Procedure:**
  - Step 1: For every  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Bor}}(\mathcal{G})$  in increasing order according to  $<_{\text{deglex}}$  compute  $\text{NF}(b(Y)^{|\alpha|} X^\alpha) = \sum_{X^\delta \in \overline{\text{Mon}}(\mathcal{G})} \lambda_{\alpha, \delta}(Y) b(Y)^{|\alpha|} X^\delta$  as follows:
    - If  $X^\alpha = X_i^{d_i}$ ,  $\text{NF}(b(Y)^{|\alpha|} X^\alpha) := -b(Y)^{d_i-1} \bar{b}_i(Y) Q_i$ .

- Else, if  $X^\alpha = X_j X^\beta$  for some  $j = 1, \dots, k$ , with  $X^\beta \in \text{Bor}(\mathcal{G})$ , define
 
$$\text{NF}(b(Y)^{|\alpha|} X^\alpha) := \sum_{(X^\gamma, X^\delta) \in \text{Mon}(\mathcal{G})^2} \lambda_{\beta, \gamma}(Y) \lambda_{\gamma', \delta}(Y) b(Y)^{|\gamma'|} X^\delta,$$

with  $X_j X^\gamma = X^{\gamma'}$ , and

$$\begin{aligned} \text{NF}(b(Y)^{|\beta|} X^\beta) &= \sum_{X^\gamma \in \text{Mon}(\mathcal{G})} \lambda_{\beta, \gamma}(Y) b(Y)^{|\gamma|} X^\gamma, \\ \text{NF}(b(Y)^{|\gamma'|} X^{\gamma'}) &= \sum_{X^\delta \in \text{Mon}(\mathcal{G})} \lambda_{\gamma', \delta}(Y) b(Y)^{|\delta|} X^\delta. \end{aligned}$$

- Step 2: Construct the matrices  $M'_1, \dots, M'_k$  corresponding to multiplication by  $b X_1, \dots, X b_k$ , expressed in the basis  $\overline{\text{Mon}}(\mathcal{G})$ , using the normal forms of elements of  $\overline{\text{Bor}}(\mathcal{G})$  already computed.
- Step 3: For every  $X^\alpha \in \text{Tab}(\mathcal{G}) \setminus (\text{Mon}(\mathcal{G}) \cup \text{Bor}(\mathcal{G}))$  in increasing order according to  $<_{\text{deglex}}$  compute  $\text{NF}(b(Y)^{|\alpha|} X^\alpha)$  as follows: since  $X^\alpha = X_j X^\beta$ , for some  $j = 1, \dots, k$ , compute the vector

$$\text{NF}(b(Y)^{|\alpha|} X^\alpha) = M'_j \cdot \text{NF}(b(Y)^{|\beta|} X^\beta).$$

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 12.9 (Special Multiplication Table). □

The complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table) uses the following lemmas.

Define

$$\begin{aligned} \overline{\text{Mon}}_{<d} &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid |\alpha| < d \right\}, \\ \overline{\text{Mon}}_{\leq d} &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid |\alpha| \leq d \right\}, \\ \overline{\text{Mon}}_d &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid |\alpha| = d \right\}, \\ \overline{\text{Mon}}_{<d}(\mathcal{G}) &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid X^\alpha \in \overline{\text{Mon}}(\mathcal{G}), |\alpha| < d \right\} \\ \overline{\text{Mon}}_{\leq d}(\mathcal{G}) &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid X^\alpha \in \overline{\text{Mon}}(\mathcal{G}), |\alpha| \leq d \right\} \\ \overline{\text{Mon}}_d(\mathcal{G}) &= \left\{ b(Y)^{|\alpha|} X^\alpha \mid X^\alpha \in \overline{\text{Mon}}(\mathcal{G}), |\alpha| = d \right\} \end{aligned}$$

**Lemma 12.11.** *The normal form of every  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_d \setminus \overline{\text{Mon}}_d(\mathcal{G})$  is a linear combination of elements of  $\overline{\text{Mon}}(\mathcal{G})_{<d}$  with coefficients in  $D[Y]$ .*

*The normal form of every  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_d$  is a linear combination of elements of  $\overline{\text{Mon}}(\mathcal{G})_{\leq d}$  with coefficients in  $D[Y]$ .*

**Proof:** We prove the result by induction on  $d$ . Suppose that the result is true for  $d' < d$ , and take  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_d \setminus \overline{\text{Mon}}_d(\mathcal{G})$

- If  $X^\alpha$  is one of the  $X_i^{d_i}$ , the result is true, since  $\text{NF}(b(Y)^{d_i} X_i^{d_i}) = -b(Y)^{d_i-1} \bar{b}_i(Y) Q_i$ .

– If  $X^\alpha$  is not one of the  $X_i^{d_i}$ , then  $X^\alpha = X_i X^\beta$  with  $b(Y)^{|\beta|} X^\beta \in \overline{\text{Mon}}_{d-1} \setminus \overline{\text{Mon}}(\mathcal{G})$ . According to the induction hypothesis, the normal form of  $b(Y)^{|\beta|} X^\beta$  is a linear combination with coefficients in  $D$  of elements  $b(Y)^{|\gamma|} X^\gamma$  of  $\overline{\text{Mon}}(\mathcal{G})_{<d-1}$ . Finally, if  $b(Y)^{|\gamma|} X^\gamma \in \overline{\text{Mon}}(\mathcal{G})_{<d-1}$ , then  $b(Y)^{|\gamma|+1} X_i X^\gamma \in \overline{\text{Mon}}_{<d}$ , and we can again use the induction hypothesis.

The last claim follows since when  $X^\alpha \in \text{Mon}(\mathcal{G})$ ,  $\text{NF}(X^\alpha) = X^\alpha$ . □

**Lemma 12.12.** *Let  $\lambda$  be a bound on the degree in  $Y$  of  $b_1, \dots, b_k$  and the coefficients of  $Q_1, \dots, Q_k$ . The entries of the matrix  $M'_i$  corresponding to multiplication by  $b(Y)X_i$  have degrees in  $Y$  at most  $k d_k (d_1 + \dots + d_{k-1} - k + 1) \lambda$ .*

**Proof:** Note that the degree in  $Y$  of  $b(Y)$  is bounded by  $k \lambda$ . For  $i = 1, \dots, k$ , let  $f_{d,i}$  be the mapping sending a polynomial  $P$  of degree  $< d$ , with coefficients in  $D[Y]$  in the basis  $\overline{\text{Mon}}_{<d}$ , to  $\text{NF}(b(Y) X_i P)$ . Note that  $\text{NF}(b(Y) X_i P)$  is a linear combination of monomials of  $\overline{\text{Mon}}_{\leq d}(\mathcal{G})$  by Lemma 12.11. We are going to estimate, for

$$d_k \leq d \leq d_1 + \dots + d_k - k,$$

the degrees in  $Y$  of the entries of the matrix  $M'_{d,i}$  of  $f_{d,i}$  expressed in the bases  $\overline{\text{Mon}}_{<d}$  and  $\overline{\text{Mon}}_{\leq d}(\mathcal{G})$ . We are going to prove by induction on  $d$  that the degrees in  $Y$  of the entries of the matrices  $M'_{d,i}$  are bounded by  $k d_1 (d - d_k + 1) \lambda$ .

If  $d = d_k$ , and  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_{<d}$ , the normal form of  $b(Y)^{|\alpha|+1} X_i X^\alpha$  is either

$$\begin{cases} -b(Y)^{d_i-1} \bar{b}_i(Y) Q_i & \text{if } d_i = d_k, X^\alpha = X_i^{d_i-1}, \\ b(Y)^{|\alpha|+1} X_i X^\alpha & \text{otherwise.} \end{cases}$$

Thus, the degrees in  $Y$  of the entries of  $N_{d_k,i}$  are bounded by  $k d_k \lambda \leq k d_1 \lambda$ .

Consider  $d_k < d < d_1 + \dots + d_k - k$  and suppose by induction hypothesis that the degrees in  $Y$  of the entries of the  $M'_{d,i}$  are bounded by

$$k d_1 (d - d_k + 1) \lambda.$$

Let  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_d \setminus \overline{\text{Mon}}(\mathcal{G})$ . Then

$$b(Y)^{|\alpha|} X^\alpha = b(Y)^{d_j} X_j^{d_j} b(Y)^{|\beta|} X^\beta$$

for some  $j = 1, \dots, k$ . Replacing  $b(Y)^{d_j} X_j^{d_j}$  by  $-b(Y)^{d_j-1} \bar{b}_j(Y) Q_j$  gives a polynomial  $R$  of total degree in  $X \leq d$  and with degree in  $Y$  bounded by  $k d_j \lambda$ . The normal form of  $b(Y)^{|\alpha|+1} X_i X^\alpha$  is the normal form of  $b(Y) X_i R$ , and is computed by multiplying the matrix  $M'_{d,i}$  with the vector of coefficients of  $R$  in the basis  $\overline{\text{Mon}}_{\leq d}$ . Thus, since  $d_j \leq d_1$ , the degrees in  $Y$  of the entries of  $M'_{d+1,i}$  are bounded by  $k d_1 ((d + 1) - d_k + 1) \lambda$ , using the complexity analysis of Algorithm 8.13 (Multiplication of matrices).



Finally we have proved that the entries of the matrix  $M'_i$  corresponding to multiplication by  $b(Y)X_i$  in  $\overline{\text{Mon}}(\mathcal{G})$ , which is a submatrix of  $M'_{d_1+\dots+d_k-k,i}$  have degrees in  $Y$  at most  $k d_1 (d_1 + \dots + d_{k-1} - k + 1) \lambda$ .  $\square$

**Lemma 12.13.** *Let  $\lambda$  be a bound on the degrees in  $Y$  of  $b_1, \dots, b_k$  and the coefficients of  $Q_1, \dots, Q_k$ ,  $\tau$  a bound on the bitsizes of  $b_1, \dots, b_k$  and the coefficients of  $Q_1, \dots, Q_k$ , and  $\nu'$  a bound on the bitsize of  $((d_1 + \dots + d_k) \lambda + 1)^\ell$ . The entries of the matrix  $M'_i$  corresponding to multiplication by  $b(Y)X_i$  in the basis  $\overline{\text{Mon}}(\mathcal{G})$  have degrees in  $Y$  at most  $k d_1 (d_1 + \dots + d_k - k) \lambda$  and bitsizes at most*

$$(d_1 + \dots + d_{k-1} - k + 1) (k d_1 + 1) (\tau + \nu').$$

**Proof:** The claim about the degrees is already proved in Lemma 12.12. The bitsize estimate is proved using the same technique. Note that, using the complexity analysis of Algorithm 8.4, the bitsizes of the coefficients of  $b$  is bounded by  $k\tau + \ell (\log_2(k \lambda) + 1)$ . For  $i = 1, \dots, k$ , let  $f_{d,i}$  be the mapping sending a polynomial  $P$ , of degree  $< d$ , with coefficients in  $D[Y]$  in the basis  $b^{|\alpha|} X^\alpha$ , to  $\text{NF}(b X_i P)$ . Note that  $\text{NF}(b X_i P)$  is a linear combination of monomials of  $\overline{\text{Mon}}_{\leq d}(\mathcal{G})$  by Lemma 12.11. We are going to estimate for  $d_k \leq d \leq d_1 + \dots + d_k - k$  the bitsizes of the coefficients the entries of the matrix  $M'_{d,i}$  of  $f_{d,i}$  expressed in the bases  $\overline{\text{Mon}}_{< d}$  and  $\overline{\text{Mon}}_{\leq d}(\mathcal{G})$ . We are going to prove by induction on  $d$  that the bitsizes of the entries of the matrices  $M'_{d,i}$  are bounded by  $(d - k + 1) (k d_1 + 1) (\tau + \nu')$ .

If  $d = d_k$ , and  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_{< d}$ , the normal form of  $b(Y)^{|\alpha|+1} X_i X^\alpha$  is

$$\begin{cases} -b(Y)^{d_i-1} \bar{b}_i(Y) Q_i & \text{if } d_i = d_k, X^\alpha = X_i^{d_i-1}, \\ b(Y)^{|\alpha|+1} X_i X^\alpha & \text{otherwise.} \end{cases}$$

Thus the bitsizes of the coefficients of the entries of  $M'_{d_k,i}$  are bounded by  $k d_k (\tau + \nu') \leq (k d_1 + 1) (\tau + \nu')$ .

Consider  $d_k < d < d_1 + \dots + d_k - k$  and suppose by induction hypothesis that the bitsizes of the entries of the matrices  $M'_{d,i}$  are bounded by

$$(d - k + 1) (k d_1 + 1) (\tau + \nu').$$

Let  $b(Y)^{|\alpha|} X^\alpha \in \overline{\text{Mon}}_d \setminus \overline{\text{Mon}}(\mathcal{G})$ , then  $b(Y)^{|\alpha|} X^\alpha = b^{d_j} X_j^{d_j} b^{|\beta|} X^\beta$  for some  $j = 1, \dots, k$ . Replacing  $b(Y)^{d_j} X_j^{d_j}$  by  $-b(Y)^{d_j-1} \bar{b}_j Q_j$  gives a polynomial  $R$  of total degree in  $X \leq d$  with coefficients of bitsizes bounded by  $k d_k (\tau + \nu')$ . The normal form of  $b(Y)^{|\alpha|+1} X_i X^\alpha$  is the normal form of  $b(Y) X_i R$ , and is computed by multiplying the matrix  $M'_{d,i}$  with the vector of coefficients of  $R$  expressed in the basis  $\overline{\text{Mon}}_{\leq d}$ . Thus, since  $d_j \leq d_1$ , the bitsizes of the entries of  $N_{d+1,i}$  are bounded by

$$((d + 1) - d_k + 1) (k d_1 + 1) (\tau + \nu'),$$

using the complexity analysis of Algorithm 8.14 (Multiplication of several matrices).

Finally we have proved that the entries of the matrix  $M'_i$  corresponding to multiplication by  $b(Y)X_i$  in  $\overline{\text{Mon}}(\mathcal{G})$ , which is a submatrix of  $M'_{d_1+\dots+d_k-k,i}$  have bitsizes at most  $(d_1 + \dots + d_{k-1} - k + 1)(k d_1 + 1)(\tau + \nu')$ .  $\square$

**Complexity analysis:**

The number of arithmetic operations in  $D[Y]$  of the algorithm is in  $O(2^k d_1 \dots d_k)^3$  using the complexity analysis of Algorithm 12.9.

If the coefficients of the polynomials in  $\mathcal{G}$  are polynomials in  $Y$  of degree bounded by  $\lambda$  we estimate the degrees in  $Y$  of the normal forms computed through the algorithm. Closely following the algorithm would give a degree in  $Y$  exponential in  $d$  since the degree looks like it is doubled each time the degree is increased by 1. So we have to proceed in a more careful way, taking into account the special structure of the Gröbner basis and using Lemma 12.12.

By Lemma 12.12, the entries of the matrix  $M'_i$  of multiplication by  $b(Y)X_i$  in  $\overline{\text{Mon}}(\mathcal{G})$ , have degree in  $Y$  at most  $k d_1 (d_1 + \dots + d_{k-1} - k + 1) \lambda$ .

Multiplying at most  $2 (d_1 + \dots + d_k - k)$  times matrices of size  $d_1 \dots d_k$  with entries of degree  $k d_1 (d_1 + \dots + d_{k-1} - k + 1) \lambda$  in  $Y$  produces matrices with entries of degree in  $Y$  bounded by  $2 k d_1 (d_1 + \dots + d_k - k)^2 \lambda$ , using the complexity analysis of Algorithm 8.13 (Multiplication of matrices).

Finally, the number of arithmetic operations in  $D$  is  $(d_1 \dots d_{k_*} \lambda)^{O(\ell)}$  since the number of arithmetic operations in  $D[Y]$  is  $(d_1 \dots d_k)^{O(1)}$ , and the degrees in  $Y$  of the polynomials appearing in the intermediate computations are bounded by  $2 k d_1 (d_1 + \dots + d_k - k)^2 \lambda$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{G}$  are bounded by  $\tau$ , multiplying at most  $2 (d_1 + \dots + d_k - k)$  times matrices of size  $d_1 \dots d_k$  with entries of degree in  $Y$   $k d_1 (d_1 + \dots + d_{k-1} - k + 1) \lambda$  and coefficients of bitsizes  $(d_1 + \dots + d_{k-1} - k + 1)(k d_1 (\tau + \nu) + \nu')$  produces matrices with entries of degree in  $Y$  bounded by  $2 k d_1 (d_1 + \dots + d_k - k)^2 \lambda$ , and coefficients of bitsizes  $2 (d_1 + \dots + d_k - k + 1)^2 (k d_1 + 1)(\tau + 4 \nu')$ , using the complexity analysis of Algorithm 8.14 (Multiplication of several matrices), since  $\nu'$  is a bound on the bitsize of  $d_1 \dots d_k$  and  $2 \nu'$  is a bound on the bitsize of  $(k d_1 (d_1 + \dots + d_k - k)^2 \lambda + 1)^\ell$ .  $\square$

## 12.4 Univariate Representation

In this section, we describe a method, based on trace computations, for solving a system of polynomial equations, in the following sense. We are going to describe the coordinates of the solutions of a zero-dimensional polynomial system as rational functions of the roots of a univariate polynomial. As before, let  $\mathcal{P} \subset K[X_1, \dots, X_k]$  be a zero-dimensional polynomial system, i.e. a finite set of polynomials such that  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  is finite. According to Theorem 4.85,  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$  is a finite dimensional vector space over  $K$  having dimension  $N \geq n = \#\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .

For any element  $a \in A$ , let  $\chi(a, T)$  be the characteristic polynomial of the linear transformation  $L_a$  from  $A$  to  $A$  defined by  $L_a(g) = ag$ . Then, according to Theorem 4.97 (Stickelberger), Equation (4.8),

$$\chi(a, T) = \prod_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} (T - a(x))^{\mu(x)}, \tag{12.2}$$

where  $\mu(x)$  is the multiplicity of  $x$  as a root of  $\chi(a, T)$ . By Remark 4.98, we have  $\chi(a, T) \in K[T]$ .

If  $a$  is separating (see Definition 4.88), for every root  $t$  of  $\chi(a, T)$ , there is a single point  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$  such that  $a(x) = t$ . Hence, it is natural to express the coordinates of the elements  $x$  in  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  as values of rational functions at the roots of  $\chi(a, T)$  when  $a$  in  $A$  is a separating element.

*Remark 12.14.* An example of this situation has already been seen in the algebraic proof of the fundamental theorem of algebra seen in Chapter 2 (proof of  $a) \Rightarrow b)$  in Theorem 2.11, see Remark 2.17). Using the notation of the proof of Theorem 2.11, the value  $z$  such that  $D(z) \neq 0$  was such that the images of all the  $\gamma_{i,j} = x_i + x_j + z x_i x_j$  where distinct and both  $x_i + x_j$  and  $x_i x_j$  where expressed as rational function of  $\gamma_{i,j}$ . □

For any  $a$  and  $f$  in  $A$ , we define

$$\varphi(a, f, T) = \sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \mu(x) f(x) \prod_{\substack{t \in \text{Zer}(\chi(a, T), \mathbb{C}) \\ t \neq a(x)}} (T - t). \tag{12.3}$$

Note that  $\chi(a, T)$  and  $\varphi(a, 1, T)$  are coprime.

If  $a$  is separating,

$$\varphi(a, f, T) = \sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \mu(x) f(x) \prod_{\substack{y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k) \\ y \neq x}} (T - a(y)),$$

and thus if  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$

$$\varphi(a, f, a(x)) = \mu(x) f(x) \prod_{\substack{y \in \text{Zer}(\mathcal{P}, \mathbb{C}^k) \\ y \neq x}} (a(x) - a(y)).$$

Hence, if  $a$  is separating,

$$\frac{\varphi(a, f, a(x))}{\varphi(a, 1, a(x))} = f(x). \tag{12.4}$$

Choosing the polynomial  $X_i$  for  $f$  in Equation (12.4), we see that

$$\frac{\varphi(a, X_i, a(x))}{\varphi(a, 1, a(x))} = x_i. \tag{12.5}$$

In other words, if  $a$  is separating and  $x$  in  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , then  $x_i$  is the value of the rational function  $\varphi(a, X_i, T) / \varphi(a, 1, T)$  at the root  $a(x)$  of the polynomial  $\chi(a, T)$ .

Note that, if  $a$  is separating, the roots of  $\mathcal{P}$  in  $C^k$  are all simple if and only if  $\chi(a, T)$  is separable, in which case

$$\varphi(a, 1, T) = \chi(a, T)'. \tag{12.6}$$

Note that for any  $a$ , not necessarily separating,

$$\frac{\varphi(a, f, a(x))}{\varphi(a, 1, a(x))} = \frac{\sum_{\substack{y \in \text{Zer}(\mathcal{P}, C) \\ a(y) = a(x)}} \mu(y) f(y)}{\sum_{\substack{y \in \text{Zer}(\mathcal{P}, C) \\ a(y) = a(x)}} \mu(y)}. \tag{12.7}$$

For any  $a$  (not necessarily separating) and  $f$  in  $A$ ,  $\varphi(a, f, T)$  belongs to  $C[T]$  by definition. In fact, as we now show, it belongs to  $K[T]$ .

**Lemma 12.15.** *If  $a$  and  $f$  are in  $A$ ,  $\varphi(a, f, T) \in K[T]$ .*

**Proof:** Since  $\chi(a, T) \in K[T]$  by Remark 4.98, let  $\bar{\chi}(a, T) \in K[T]$  be the monic separable part of  $\chi(a, T)$ . It is clear that  $\bar{\chi}(a, T) \in K[T]$ . Then

$$\begin{aligned} \frac{\varphi(a, f, T)}{\bar{\chi}(a, T)} &= \sum_{x \in \text{Zer}(\mathcal{P}, C^k)} \frac{\mu(x) f(x)}{T - a(x)}, \\ \bar{\chi}(a, T) &= \prod_{t \in Z(\chi(a, T), C)} (T - t), \\ \frac{\varphi(a, f, T)}{\bar{\chi}(a, T)} &= \prod_{t \in Z(\chi(a, T), C)} (T - t), \\ &= \sum_{i \geq 0} \sum_{x \in \text{Zer}(\mathcal{P}, C^k)} \frac{\mu(x) f(x) a(x)^i}{T^{i+1}} \\ &= \sum_{i \geq 0} \frac{\text{Tr}(L_{fa^i})}{T^{i+1}}. \end{aligned}$$

Let

$$\bar{\chi}(a, T) = \sum_{j=0}^{n'} c_{n'-j} T^{n'-j},$$

with  $c_{n'-j} \in K$ ,  $c_{n'} = 1$ ,  $n' \leq n$ . Note that if  $a$  is separating, then  $n' = n$ . Multiplying both sides by  $\bar{\chi}(a, T)$ , which is in  $K[T]$  and using the fact that  $\varphi(a, f, T)$  is a polynomial in  $C[T]$ , we have

$$\varphi(a, f, T) = \sum_{i=0}^{n'-1} \sum_{j=0}^{n'-i-1} \text{Tr}(L_{fa^i}) c_{n'-j} T^{n'-i-1-j} \tag{12.8}$$

Consider a Gröbner basis  $\mathcal{G}$  of  $\text{Ideal}(\mathcal{P}, K)$  and express  $L_{fa^i}$  in the basis  $\text{Mon}(\mathcal{G})$ . Then  $\text{Tr}(L_{fa^i})$ , which is the trace of a matrix with entries in  $K$ , is in  $K$ . This proves  $\varphi(a, f, T) \in K[T]$ . □

The previous discussion suggests the following definition and proposition.

A  **$k$ -univariate representation**  $u$  is a  $k + 2$ -tuple of polynomials in  $K[T]$ ,

$$u = (f(T), g(T)), \text{ with } g = (g_0(T), g_1(T), \dots, g_k(T)),$$

such that  $f$  and  $g_0$  are coprime. Note that  $g_0(t) \neq 0$  if  $t \in C$  is a root of  $f(T)$ . The **points associated** to a univariate representation  $u$  are the points

$$x_u(t) = \left( \frac{g_1(t)}{g_0(t)}, \dots, \frac{g_k(t)}{g_0(t)} \right) \in C^k \tag{12.9}$$

where  $t \in C$  is a root of  $f(T)$ .

Let  $\mathcal{P} \subset K[X_1, \dots, X_k]$  be a finite set of polynomials such that  $\text{Zer}(\mathcal{P}, C^k)$  is finite. The  $k + 2$ -tuple  $u = (f(T), g(T))$ , **represents**  $\text{Zer}(\mathcal{P}, C^k)$  if  $u$  is a univariate representation and

$$\text{Zer}(\mathcal{P}, C^k) = \{x \in C^k \mid \exists t \in \text{Zer}(f, C) \ x = x_u(t)\}.$$

A **real  $k$ -univariate representation** is a pair  $u, \sigma$  where  $u$  is a  $k$ -univariate representation and  $\sigma$  is the Thom encoding of a root of  $f$ ,  $t_\sigma \in R$ . The **point associated** to the real univariate representation  $u, \sigma$  is the point

$$x_u(t_\sigma) = \left( \frac{g_1(t_\sigma)}{g_0(t_\sigma)}, \dots, \frac{g_k(t_\sigma)}{g_0(t_\sigma)} \right) \in R^k. \tag{12.10}$$

Let

$$\varphi(a, T) = (\varphi(a, 1, T), \varphi(a, X_1, T), \dots, \varphi(a, X_k, T)) \tag{12.11}$$

**Proposition 12.16.** *Let  $\mathcal{P} \subset K[X_1, \dots, X_k]$  a zero-dimensional system and  $a \in A$ .*

*The  $k$ -univariate representation  $(\chi(a, T), \varphi(a, T))$  represents  $\text{Zer}(\mathcal{P}, C^k)$  if and only if  $a$  is separating.*

*If  $a$  is separating, the following properties hold.*

- *The degree of the separable part of  $\chi(a, T)$  is equal to the number of elements in  $\text{Zer}(\mathcal{P}, C^k)$ .*
- *The bijection  $x \mapsto a(x)$  from  $\text{Zer}(\mathcal{P}, C^k)$  to  $\text{Zer}(\chi(a, T), R)$  respects the multiplicities.*

**Proof:** If  $a$  is separating,  $(\chi(a, T), \varphi(a, T))$  represents  $\text{Zer}(\mathcal{P}, C^k)$  by (12.5). Conversely if  $(\chi(a, T), \varphi(a, T))$  represents  $\text{Zer}(\mathcal{P}, C^k)$ , then  $a(x) = a(y)$  imply  $x = y$ , hence  $a$  is separating.

Since the degree of the separable part of  $\chi(a, T)$  is equal to the number of distinct roots of  $\chi(a, T)$ , it coincides with  $\text{Zer}(\mathcal{P}, C^k)$  when  $a$  is separating.

Finally, if  $a$  is separating, the multiplicity of  $a(x)$  is equal to  $\mu(x)$  by (12.2). □

The following proposition gives a useful criterion for  $a$  to be separating.

**Proposition 12.17.** *Let*

*The following properties are equivalent:*

- a) *The element  $a \in A$  is separating.*
- b) *The  $k + 2$ -tuple  $(\chi(a, T), \varphi(a, T))$  represents  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .*
- c) *For every  $k = 1, \dots, k$ ,*

$$\varphi(a, 1, a) X_i - \varphi(a, X_i, a) \in \sqrt{\text{Ideal}(\mathcal{P}, \mathbb{K})}.$$

**Proof:** By Proposition 12.16, a) and b) are equivalent.

Let us prove that b) and c) are equivalent. Since  $\chi(a, T)$  and  $\varphi(a, 1, T)$  are coprime, for every  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , and every  $i = 1, \dots, k$

$$x_i = \frac{\varphi(a, X_i, a(x))}{\varphi(a, 1, a(x))}$$

is equivalent to the property (P):  $\varphi(a, 1, a) X_i - \varphi(a, X_i, a)$  vanishes at every  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ , for every  $i = 1, \dots, k$ . Property (P) is equivalent to

$$\varphi(a, 1, a) X_i - \varphi(a, X_i, a) \in \sqrt{\text{Ideal}(\mathcal{P}, \mathbb{K})}$$

for every  $i = 1, \dots, k$  by Hilbert Nullstellensatz (Theorem 4.78). □

Since  $a \in A = \mathbb{K}[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, \mathbb{K})$  we also obtain,

**Corollary 12.18.**

*If  $a \in A$  is separating,*

$$a(\text{Zer}(\mathcal{P}, \mathbb{R}^k)) = \text{Zer}(\chi(a, T), \mathbb{R}).$$

*In particular,  $\#\text{Zer}(\mathcal{P}, \mathbb{R}^k) = \#\text{Zer}(\chi(a, T), \mathbb{R})$ .*

**Proof:** Since  $a \in A$ , if  $x \in \text{Zer}(\mathcal{P}, \mathbb{R}^k)$ ,  $a(x) \in \text{Zer}(\chi(a, T), \mathbb{R})$ . Conversely, if  $\chi(a, t) = 0$ , then  $t = a(x)$  for  $x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)$ ,  $x = x_{u(\mathcal{P}, a)}(t)$  and

$$x_i = \frac{\varphi(a, X_i, a(x))}{\varphi(a, 1, a(x))},$$

by (12.5). Since  $\varphi(a, X_i, T)$  and  $\varphi(a, 1, T)$  belong to  $\mathbb{K}[T]$  by Lemma 12.15,  $t = a(x) \in \mathbb{R}$  implies  $x \in \mathbb{R}^k$ . □

Let  $D$  be a ring contained in  $\mathbb{K}$ ,  $\mathcal{P} \subset D[X_1, \dots, X_k]$  a zero-dimensional system and  $\mathcal{B}$  a basis of  $A$  such that the multiplication table of  $A$  in  $\mathcal{B}$  has entries in  $D$ . Consider  $a \in A$ ,  $b \in D$ , and suppose that  $a, b, bX_1, \dots, bX_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ . Note that  $\chi(a, T) \in D[T]$ , and let  $\bar{\chi}(a, T)$  be a separable part of  $\chi(a, T)$  with coefficients in  $D$  and leading coefficient  $c$ . We denote

$$\varphi_b(a, T) = (\varphi(a, b, T), \varphi(a, bX_1, T), \dots, \varphi(a, bX_k, T)), \tag{12.12}$$

(using (12.2) and (?)).

Since

$$\frac{c \varphi(a, b X_i, a(x))}{c \varphi(a, b, a(x))} = \frac{\varphi(a, X_i, a(x))}{\varphi(a, 1, a(x))} = x_i. \tag{12.13}$$

Proposition 12.16, Proposition 12.17 and Corollary 12.18 hold when replacing  $\varphi(a, T)$  by  $c \varphi_b(a, T)$ . Introducing  $c, b$  plays a role in guaranteeing that the computations take place inside  $D$ .

**Proposition 12.19.** *Let  $D$  be a ring contained in  $K$ ,  $\mathcal{P} \subset D[X_1, \dots, X_k]$  a zero-dimensional system. Consider  $a \in A$ ,  $b \neq 0 \in D$ .*

*Let  $\mathcal{B}$  a basis of  $A$  such that the multiplication table of  $A$  in  $\mathcal{B}$  has entries in  $D$ . Suppose that  $a, b, b X_1, \dots, b X_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ , then, denoting by  $c$  the leading coefficient of a separable part  $\bar{\chi}$  of  $\chi$  in  $D[T]$ , the components of  $c \varphi_b(a, T) \in D[T]^{k+1}$ .*

**Proof:** The polynomial  $\bar{\chi}(a, T)$  has coefficients in  $D$  and the various  $\text{Tr}(L_{a^i})$   $\text{Tr}(L_{b X_i a^j})$  belong to  $D$ . Thus by Equation (12.8),  $c \varphi_b(a, T) \in D[T]^{k+1}$ .  $\square$

Let  $a \in A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ . The polynomial  $\chi(a, T)$  is related to the traces of the powers of  $a$  as follows: The  $i$ -th Newton sum  $N_i$  associated to the polynomial  $\chi(a, T)$  is the sum of the  $i$ -th powers of the roots of  $\chi(a, T)$  and is thus

$$N_i = \sum_{x \in \text{Zer}(\mathcal{P}, \mathbb{C}^k)} \mu(x) a(x)^i.$$

According to Proposition 4.54,  $N_i = \text{Tr}(L_{a^i})$ . Let

$$\chi(a, T) = \sum_{i=0}^N b_{N-i} T^{N-i}.$$

According to Newton’s formula (Equation (4.1)),

$$(N - i) b_{N-i} = \sum_{j=0}^i \text{Tr}(L_{a^j}) b_{N-i+j}, \tag{12.14}$$

so that  $\chi(a, T)$  can be computed from  $\text{Tr}(L_{a^i})$ , for  $i = 0, \dots, N$ . Moreover,  $a$  is separating when the number of distinct roots of  $\chi(a, T)$  is  $n = \#\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ .

We then compute a separable part of  $\chi(a, T)$

$$\bar{\chi}(a, T) = \sum_{j=0}^{n'} c_{n'-j} T^{n'-j},$$

with leading coefficient  $c = c_{n'}$  and write  $c \varphi(a, f, T)$  as

$$c \varphi(a, f, T) = \sum_{i=0}^{n'-1} \text{Tr}(L_{f a^i}) \text{Hor}_{n'-i-1}(\bar{\chi}(a, T), T), \tag{12.15}$$

where  $\text{Hor}_i(P, T)$  is the  $i$ -th Horner polynomial associated to  $P$  (see Notation 8.6).

So,  $(\chi(a, T), c \varphi_b(a, T))$ , can easily be obtained from the following traces

$$\text{Tr}(a^i), i = 0, \dots, N, \text{Tr}(a^i b X_j), i = 0, \dots, n', j = 1, \dots, k,$$

where  $n' = \#\text{Zer}(\chi(a, T), C) \leq n = \#\text{Zer}(\mathcal{P}, C^k)$ , using Equation (12.14) and Equation (12.15).

*Algorithm 12.11. [Candidate Univariate Representation]*

- **Structure:** a ring  $D$  with division in  $\mathbb{Z}$  contained in a field  $K$ .
- **Input:** a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$  such that the multiplication table  $\mathcal{M}$  of  $A$  in  $\mathcal{B}$  has entries in  $D$ , an element  $a \in A$  and  $b \neq 0 \in D$  such that  $a, b, b X_1, \dots, b X_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ .
- **Output:**  $c \in D$ , and  $(\chi(a, T), c \varphi_b(a, T)) \in D[T^{k+2}]$ .
- **Complexity:**  $O(N^3 + k N^2)$ , where  $N$  is the number of elements of  $\mathcal{B}$ , in the special case when  $\mathcal{B} = \text{Mon}(\mathcal{G})$ .
- **Procedure:**
  - Step 1: Compute the traces of  $L_{a^i}$ ,  $i = 1, \dots, N$  using Algorithm 12.5 (Trace). Then compute the coefficients of  $\chi(a, T)$  using Algorithm 8.11 (Newton sums) and Equation (12.14).
  - Step 2: Compute the separable part of  $\chi(a, T)$  using Algorithm 10.1 (Gcd and Gcd-free Part),  $c$  its leading coefficient.
  - Step 3: Compute  $c \varphi_b(a, T)$  using Algorithm 12.5 (Trace) and Equation (12.15).
  - Return  $(\chi(a, T), c \varphi_b(a, b))$ .

**Proof of correctness:** Immediate. Note that we know in advance that  $\chi(a, T) \in D[T]$  by Corollary 12.15, so exact division by an integer is possible.  $\square$

**Complexity analysis:** Let  $N$  be the dimension of the  $K$ -vector space  $A$ .

Before computing the traces, it is necessary to compute the normal forms of  $1, a, \dots, a^n$ , which is done by multiplying at most  $N$  times the matrix of multiplication by  $a$  (which is a linear combination of the matrices  $M_i$  of multiplication by  $X_i$ ) which takes  $O(N^3)$  arithmetic operations. According to the complexity analyses of Algorithm 12.5 (Trace), Algorithm 8.11 (Newton sums), and Algorithm 10.1 (Gcd and Gcd-free part), using Equation (12.15),  $c \varphi_b(a, T)$  can clearly be computed in  $O(k N^2)$  arithmetic operations.  $\square$

*Algorithm 12.12. [Univariate Representation]*

- **Structure:** a ring  $D$  with division in  $\mathbb{Z}$  contained in a field  $K$ .
- **Input:** a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$ , such that the multiplication table  $\mathcal{M}$  of  $A$  in  $\mathcal{B}$  has entries in  $D$ , and  $b \neq 0 \in D$  such that  $b, b X_1, \dots, b X_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ .
- **Output:** a univariate representation  $u$  representing  $\text{Zer}(\mathcal{P}, C^k)$ .
- **Complexity:**  $O(k N^2(N^3 + k N^2))$ , where  $N$  is the number of elements of  $\mathcal{B}$ , in the special case when  $\mathcal{B} = \text{Mon}(\mathcal{G})$ .



• **Procedure:**

- Compute  $n = \#Zer(\mathcal{P}, \mathbb{C}^k)$  using Algorithm 12.6 (Number of distinct zeros).
- Initialize  $i := 0$ .
- ( $\star$ ) Take  $a := X_1 + i X_2 + i^2 X_3 + \dots + i^{k-1} X_k$ . Compute  $\chi(b a, T)$  using Step 1 of Algorithm 12.11 (Candidate Univariate Representation). Compute

$$n(b a) = \deg(\chi(b a, T)) - \deg(\gcd(\chi(b a, T), \chi'(b a, T)))$$

using Algorithm 8.21 (Signed Subresultant).

- While  $n(b a) \neq n$ ,  $i := i + 1$ , return to ( $\star$ ).
- Compute  $c$  and  $c \varphi_b(b a, T)$  using Step 3 of Algorithm 12.11 (Candidate Univariate Representation).
- Return  $u = (\chi(b a, T), c \varphi_b(b a, T))$ .

**Proof of correctness:** Let  $N$  be the dimension of the  $K$ -vector space  $A$ . We know by Lemma 4.89 that there exists a separating element

$$a := X_1 + i X_2 + i^2 X_3 + \dots + i^{k-1} X_k$$

for  $i \leq (k - 1) \binom{n}{2}$ , and by Theorem 4.85 that  $n \leq N$ . Note that  $b a$  is separating as well since  $b \neq 0$ . The number  $n(b a)$  is the number of distinct roots of  $\chi(b a, T)$ , and  $n(b a) = n$  if and only if  $a$  is separating.  $\square$

**Complexity analysis:** The number of different  $a$  to consider is

$$(k - 1) \binom{N}{2} + 1,$$

and for each  $a$  the cost of computation is  $O(N^3 + k N^2)$  according to the complexity analysis of Algorithm 12.11 (Candidate Univariate Representation). Thus the complexity is  $O(k N^2(N^3 + k N^2))$ .  $\square$

*Remark 12.20.* Algorithm 12.12 (Univariate Representation) can be improved in various ways [134]. In particular, rather than looking for separating elements which are linear combination of variables, it is preferable to check first whether variables are separating. Second, the computations of Algorithm 12.6 (Number of Distinct Zeros) as well as the computation of a separating element can be performed using modular arithmetic, which avoids any growth of coefficients. Of course, if the prime modulo which the computations is performed is unlucky, the number of distinct elements and the separating element are not computed correctly. So it is useful to have a test that a candidate separating element is indeed separating; this can be checked quickly using Proposition 12.17 and Theorem 4.99. Similar remarks apply to Algorithm 12.13. Efficient computations of univariate representations can be found in [135].  $\square$

When we know in advance that the zeroes of the polynomial system are all simple, which will be the case in many algorithms in the next chapters, the computation above can be simplified.

*Algorithm 12.13.* [Simple Univariate Representation]

- **Structure:** a ring  $D$  with division in  $\mathbb{Z}$  contained in a field  $K$ .
- **Input:** a zero-dimensional polynomial system  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a basis  $\mathcal{B}$  of  $A = K[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K)$  such that the multiplication table  $\mathcal{M}$  of  $A$  in  $\mathcal{B}$  has entries in  $D$ , and  $b \neq 0 \in D$  such that  $b, bX_1, \dots, bX_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ . Moreover all the zeros of  $\text{Zer}(\mathcal{P}, C^k)$  are simple, so that  $N = n$ .
- **Output:** a univariate representation  $u$  representing  $\text{Zer}(\mathcal{P}, C^k)$ .
- **Complexity:**  $O(kN^2(N^3 + kN^2))$ , where  $N$  is the number of elements of  $\mathcal{B}$ , in the special case when  $\mathcal{B} = \text{Mon}(\mathcal{G})$ .
- **Procedure:**
  - Initialize  $i := 0$ .
  - ( $\star$ ) Take  $a := X_1 + iX_2 + i^2X_3 + \dots + i^{k-1}X_k$ . Compute  $\chi(a, T)$  using Step 1 of Algorithm 12.11 (Candidate Univariate Representation).
  - Compute  $\text{gcd}(\chi(ba, T), \chi'(ba, T))$  by Algorithm 8.21 (Signed Subresultant).
  - While  $\text{deg}(\text{gcd}(\chi(ba, T), \chi'(ba, T))) \neq 0$ ,  $i := i + 1$ , return to ( $\star$ ).
  - Compute  $\varphi_b(ba, b, T)$  by Step 3 of Algorithm 12.11 (Candidate Univariate Representation).
  - Return  $u = (\chi(ba, T), \varphi_b(ba, b, T))$ .

**Proof of correctness:** Let  $N$  be the dimension of the  $K$ -vector space  $A$ . Since all the zeros of  $\text{Zer}(\mathcal{P}, C^k)$  are simple,  $n = N$  by Theorem 4.85. We know by Lemma 4.89 that there exists a separating element  $i \leq (k-1) \binom{N}{2}$ . Since all the zeros of  $\text{Zer}(\mathcal{P}, C^k)$  are simple, and  $b \neq 0$ ,  $a$  is separating if and only if  $\chi(ba, T)$  is separable.  $\square$

**Complexity analysis:** The number of different  $a$  to consider is

$$(k-1) \binom{N}{2} + 1,$$

and for each  $a$  the computation is  $O(N^3 + kN^2)$  according to the complexity analysis of Algorithm 12.11 (Univariate Representation). Thus the complexity is  $O(kN^2(N^3 + kN^2))$ .  $\square$

*Remark 12.21.* It is clear what we can use the rational univariate representation to give an alternative method to Algorithm 12.7 (Multivariate Tarski-query) (and to Algorithm 12.8 (Multivariate Sign Determination)) in the multivariate case. Given a univariate representation  $(f, g)$  representing  $\text{Zer}(\mathcal{P}, C^k)$  we simply replace the Tarski-query  $\text{TaQ}(Q, \text{Zer}(\mathcal{P}, R^k))$  by the Tarski-query  $\text{TaQ}(Q_u, \text{Zer}(f, R))$ , with

$$Q_u = g_0^e Q \left( \frac{g_k}{g_0}, \dots, \frac{g_k}{g_0} \right). \quad (12.16) \square$$

*Remark 12.22.* Note that in the two last sections, the computation of the traces of various  $L_f$  was crucial in most algorithms. These are easy to compute once a multiplication table is known. However, in big examples, the size of the multiplication table can be the limiting factor in the computations. Efficient ways for computing the traces without storing the whole multiplication table are explained in [134, 136].  $\square$

## 12.5 Limits of the Solutions of a Polynomial System

In the next chapters, it will be helpful for complexity reasons to perturb polynomials, making infinitesimal deformations. The solutions to systems of perturbed equations belong to fields of algebraic Puiseux series. We will have to deal with the following problem: given a finite set of points with algebraic Puiseux series coordinates, compute the limits of the points as the infinitesimal quantities tend to zero.

**Notation 12.23. [Limit]** Let  $\varepsilon = \varepsilon_1, \dots, \varepsilon_m$  be variables. As usual, we denote by  $K[\varepsilon] = K[\varepsilon_1, \dots, \varepsilon_m]$  the ring of polynomials in  $\varepsilon_1, \dots, \varepsilon_m$ , and by  $K(\varepsilon) = K(\varepsilon_1, \dots, \varepsilon_m)$ , the field of rational functions in  $\varepsilon_1, \dots, \varepsilon_m$ , which is the fraction field of  $K[\varepsilon]$ . If  $K$  is a field of characteristic 0, and  $\delta$  is a variable we denote as in Chapter 2 by  $K\langle\delta\rangle$  the field of algebraic Puiseux series in  $\delta$  with coefficients in  $K$ . We denote by  $K\langle\varepsilon\rangle$  the field  $K\langle\varepsilon_1\rangle\dots\langle\varepsilon_m\rangle$ . If  $\nu = (\nu_1, \dots, \nu_m) \in \mathbb{Q}^m$ ,  $\varepsilon^\nu$  denotes  $\varepsilon_1^{\nu_1}\dots\varepsilon_m^{\nu_m}$ . It follows from Theorem 2.91 and Theorem 2.92 that  $R\langle\varepsilon\rangle$  is real closed and  $C\langle\varepsilon\rangle$  is algebraically closed. Note that in  $R\langle\varepsilon\rangle$ ,  $\varepsilon^\nu < \varepsilon^\mu$  if and only if  $(\nu_m, \dots, \nu_1) >_{\text{lex}} (\mu_m, \dots, \mu_1)$  (see Definition 2.14). In particular  $\varepsilon_m < \dots < \varepsilon_1$  in  $R\langle\varepsilon\rangle$ . The preceding order on elements of  $\mathbb{Q}^m$  and their corresponding monomials is denoted by  $<_\varepsilon$  to avoid confusions. We denote by  $K(\varepsilon)_b$  and  $K\langle\varepsilon\rangle_b$  the subrings of  $K(\varepsilon)$  and  $K\langle\varepsilon\rangle$  which are sums of  $\varepsilon^\nu$  with  $\varepsilon^\nu \leq_\varepsilon 1$ .

The elements of  $R(\varepsilon)_b$  and  $R\langle\varepsilon\rangle_b$  (resp.  $C(\varepsilon)_b$  and  $C\langle\varepsilon\rangle_b$ ) are the elements of  $R(\varepsilon)$  and  $R\langle\varepsilon\rangle$  (resp.  $C(\varepsilon)$  and  $C\langle\varepsilon\rangle$ ) bounded over  $R$  i.e. whose absolute value (resp. norm) is bounded by a positive element of  $R$ .

An element  $\tau \neq 0$  of  $K\langle\varepsilon\rangle$  can be written uniquely as  $\varepsilon^{o(\tau)} (\text{In}(\tau) + \tau')$  with  $\varepsilon^{o(\tau)}$  the biggest monomial of  $\tau$  for the order of  $<_\varepsilon$ ,  $\text{In}(\tau) \neq 0 \in K$  and  $\tau' \in K\langle\varepsilon\rangle_b$  with biggest monomial  $<_\varepsilon 1$ . The  $m$ -tuple  $o(\tau)$  is the **order** of  $\tau$  and  $\text{In}(\tau)$  is its initial coefficient. We have

$$\begin{aligned} o(\tau\tau') &= o(\tau) + o(\tau'), \\ o(\tau) <_\varepsilon o(\tau') &\Rightarrow o(\tau + \tau') = o(\tau'), \\ o(\tau) = o(\tau') &\Rightarrow o(\tau) \leq_\varepsilon o(\tau + \tau'). \end{aligned}$$

We define  $\lim_\varepsilon(\tau)$  from  $K\langle\varepsilon\rangle_b$  to  $K$  as follows:

$$\begin{cases} \lim_\varepsilon(\tau) = \text{In}(\tau) & \text{if } \varepsilon^{o(\tau)} = 1, \\ \lim_\varepsilon(\tau) = 0 & \text{otherwise.} \end{cases} \quad \square$$

*Example 12.24.* Let  $m = 2$  and consider  $K\langle\varepsilon_1, \varepsilon_2\rangle$ . Note that  $\varepsilon_1/\varepsilon_2 \notin K\langle\varepsilon_1, \varepsilon_2\rangle_b$ , and  $\varepsilon_2/\varepsilon_1 \in K\langle\varepsilon_1, \varepsilon_2\rangle_b$ . Then  $\lim_\varepsilon (\varepsilon_2/\varepsilon_1 + 2\varepsilon_1) = 0$ , and  $\lim_\varepsilon (\varepsilon_2/\varepsilon_1 + 2) = 2$ , while  $\lim_\varepsilon$  is not defined for  $\varepsilon_1/\varepsilon_2 \notin K\langle\varepsilon_1, \varepsilon_2\rangle_b$ .  $\square$

We first discuss how to find the limits of the roots of a univariate monic polynomial  $F(T) \in K(\varepsilon)[T]$ . Note that in our computations, we are going to compute polynomials in  $K(\varepsilon)[T]$ , with roots in  $C\langle\varepsilon\rangle[T]$ .

We denote  $\text{Zer}_b(F(T), C\langle\varepsilon\rangle)$  the set  $\{\tau \in C\langle\varepsilon\rangle_b \mid F(\tau) = 0\}$ .

**Notation 12.25. [Order of a polynomial]** Given  $F(T) \in K(\varepsilon)[T]$ , we denote by  $o(F)$  the maximal value of  $o(c)$  with respect to the ordering  $<_\varepsilon$  for  $c$  coefficient of  $F$ . In other words,  $\varepsilon^{o(F)}$  is the minimal monomial with respect to the ordering  $<_\varepsilon$  such that  $\varepsilon^{-o(F)} F(T)$  belongs to  $K(\varepsilon)_b[T]$ .  $\square$

Denote by  $f(T) = \lim_\varepsilon (\varepsilon^{-o(F)} F(T))$  the univariate polynomial obtained by replacing the coefficients of  $\varepsilon^{-o(F)} F(T)$  by their limit under  $\lim_\varepsilon$ .

Now we relate the roots of  $F(T)$  in  $C\langle\varepsilon\rangle$  and the roots of  $f(T)$  in  $C$ .

**Lemma 12.26.** *Let  $Z_b$  be the set of roots of  $F(T)$  in  $C\langle\varepsilon\rangle_b$ , and let  $Z_u$  be the roots of  $F(T)$  in  $C\langle\varepsilon\rangle \setminus C\langle\varepsilon\rangle_b$ . We denote by  $\mu(\tau)$  the multiplicity of a root  $\tau$  of*

a) *We have*

$$o(F) = \sum_{\tau \in Z_u}^p \mu(\tau) o(\tau),$$

$$\text{Zer}(f(T), C) = \lim_\varepsilon (\text{Zer}_b(F(T), C\langle\varepsilon\rangle)).$$

b) *If  $t$  is a root of multiplicity  $\mu$  of  $f(T)$  in  $C$ ,*

$$\mu = \sum_{\substack{\tau \in Z_b \\ \lim_\varepsilon(\tau) = t}} \mu(\tau).$$

**Proof:** a) We have

$$F(T) = \prod_{\tau \in Z_b} (T - \tau) \prod_{\tau \in Z_u} (T - \tau) \in K(\varepsilon)[T],$$

with  $o(\tau) \leq_\varepsilon 0$ , for  $\tau \in Z_b$ ,  $o(\tau) >_\varepsilon 0$ , for  $\tau \in Z_u$ . Using the properties of the order listed in Notation 12.23, and denoting  $\ell = \sum_{\tau \in Z_b} \mu(\tau)$ , the order of the coefficient of  $T^\ell$  in  $F(T)$  is exactly  $\sum_{\tau \in Z_u} \mu(\tau) o(\tau)$ . Moreover, the order of any other coefficient of  $F(T)$  is at most  $\sum_{\tau \in Z_u} \mu(\tau) o(\tau)$  for  $<_\varepsilon$ . Thus  $o(F) = \sum_{\tau \in Z_u}^p \mu(\tau) o(\tau)$  and

$$\varepsilon^{-o(F)} F(T) = \prod_{\tau \in Z_b} (T - \tau) \prod_{\tau \in Z_u} (\varepsilon^{-o(\tau)} T - \varepsilon^{-o(\tau)} \tau) \in K(\varepsilon)_b[T].$$

Taking  $\lim_\varepsilon$  on both sides, we get

$$f(T) = \prod_{\tau \in Z_u} (-\ln(\tau)) \prod_{\tau \in Z_u} (T - \lim_\varepsilon(\tau)).$$

b) is an immediate consequence of the last equality. □

**Corollary 12.27.** *If  $F(T)$  is separable, the number  $\deg_T(F(T)) - \deg_T(f(T))$  is the number of unbounded roots of  $P$ .*

Now, let  $\mathcal{P} \subset K(\varepsilon)[X_1, \dots, X_k]$  be a zero-dimensional polynomial system, so that  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  is non-empty and finite, and

$$A = K(\varepsilon)[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, K(\varepsilon)).$$

Suppose, moreover, for the rest of this section that all the zeros of  $\mathcal{P}$  are simple. This assumption leads to technical simplifications and will be satisfied whenever we apply the results of this section in the future. We define

$$\begin{aligned} \text{Zer}_b(\mathcal{P}, R\langle\varepsilon\rangle^k) &= \text{Zer}(\mathcal{P}, R\langle\varepsilon\rangle^k) \cap R\langle\varepsilon\rangle_b^k, \\ \text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k) &= \text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k) \cap C\langle\varepsilon\rangle_b^k. \end{aligned}$$

These are the points of  $\text{Zer}(\mathcal{P}, R\langle\varepsilon\rangle^k)$  and  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  that are bounded over  $R$ . Note that  $\lim_\varepsilon(\text{Zer}_b(\mathcal{P}, R\langle\varepsilon\rangle^k)) \subset \lim_\varepsilon(\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)) \cap R^k$ . Observe that this inclusion might be strict, since there may be algebraic Puiseux series with complex coefficients and real  $\lim_\varepsilon$ . If  $a \in K[X_1, \dots, X_k]$ ,  $a$  defines a mapping from  $\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)$  to  $C\langle\varepsilon\rangle_b$ , also denoted by  $a$ , associating to  $x$  the element  $a(x)$ .

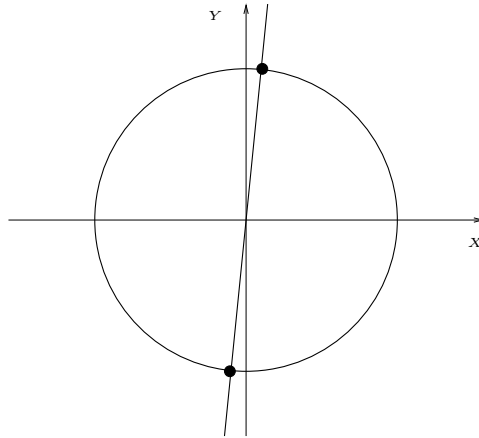
We are going to describe  $\lim_\varepsilon(\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k))$  by using a univariate representation of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  and taking its limit. In order to give such a description of  $\lim_\varepsilon(\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k))$ , it is useful to define the notion of well-separating element.

A **well-separating element**  $a$  is an element of  $K[X_1, \dots, X_k]$  that is a separating element for  $\mathcal{P}$ , such that  $a$  sends unbounded elements of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  to unbounded elements of  $C\langle\varepsilon\rangle$ , and such that  $a$  sends two non-infinitesimally close elements of  $\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)$  on two non-infinitesimally close elements of  $C\langle\varepsilon\rangle_b$ .

To illustrate how the notions of separating element and well-separating element can differ, consider the following examples:

*Example 12.28.*

- a) Consider the polynomial system  $XY = 1, X = \varepsilon$ . The only solution is  $(\varepsilon, 1/\varepsilon)$  which is unbounded. The image of this solution by  $X$  is  $\varepsilon$ , which is bounded. Thus  $X$  is separating, but not well-separating.
- b) Consider the polynomial system  $X^2 + Y^2 - 1 = 0, \varepsilon Y = X$ . The only solutions are  $(\varepsilon/(1 + \varepsilon^2)^{1/2}, 1/(1 + \varepsilon^2)^{1/2}), (-\varepsilon/(1 + \varepsilon^2)^{1/2}, -1/(1 + \varepsilon^2)^{1/2})$  which are bounded and not infinitesimally close (see Figure 12.2).



**Fig. 12.2.** Separating but not well-separating

The image of these solutions by  $X$  are  $\varepsilon/(1 + \varepsilon^2)^{1/2}, -\varepsilon/(1 + \varepsilon^2)^{1/2}$ , which are infinitesimally close. Thus  $X$  is separating, but not well-separating.  $\square$

Let  $a \in K[X_1, \dots, X_k]$  be a separating element for  $\mathcal{P}$ . Since the polynomial system  $\mathcal{P}$  is contained in  $K(\varepsilon)[X_1, \dots, X_k]$ , the polynomials of the univariate representation  $(\chi(a, T), \varphi(a, T))$  are elements of  $K(\varepsilon)[T]$ . Note that  $\chi(a, T)$  is monic and separable, since we have supposed that all the zeros of  $\mathcal{P}$  in  $C\langle \varepsilon \rangle^k$  are simple. Using Notation 12.25, note that  $\varepsilon^{-o(\chi(a, T))} \varphi(a, 1, T) \in K(\varepsilon)_b[T]$  since it is the derivative of  $\varepsilon^{-o(\chi(a, T))} \chi(a, T) \in K(\varepsilon)_b[T]$ , by Equation (12.6). However it may happen that some  $\varepsilon^{-o(\chi(a, T))} \varphi(a, X_i, T)$  do not belong to  $K(\varepsilon)_b[T]$ . In other words, denoting by  $o(\varphi(a, T))$  the maximum value for  $<_\varepsilon$  of  $o(c)$  for  $c$  a coefficient of  $\chi(a, T), \varphi(a, X_i, T), i = 1, \dots, k$ , it may happen that  $\varepsilon^{o(\varphi(a, T))} >_\varepsilon \varepsilon^{o(\chi(a, T))}$ .

*Example 12.29.* In Example 12.28 a), with  $a = X$ ,

$$\varphi(a, 1, T) = 1, \varphi(a, X_1, T) = \varepsilon, \varphi(a, X_2, T) = 1/\varepsilon.$$

Thus  $o(\chi(a, T)) = 0, o(\varphi(a, T)) = -1$ , and  $\varphi(a, X_2, T) \notin K(\varepsilon)_b[T]$ .  $\square$

**Notation 12.30.** When  $o(\chi(a, T)) = o(\varphi(a, T))$  denote by  $(\hat{\chi}(a, T), \hat{\varphi}(a, T))$  the  $j + 2$ -tuple defined by

$$\begin{aligned} \hat{\chi}(a, T) &= \lim_{\varepsilon} (\varepsilon^{-o(\chi(a, T))} \chi(a, T)), \\ \hat{\varphi}(a, T) &= \lim_{\varepsilon} (\varepsilon^{-o(\chi(a, T))} \varphi(a, T)), \end{aligned}$$

with  $\hat{\varphi}(a, T) = (\hat{\varphi}(a, 1, t), \hat{\varphi}(a, X_1, t), \dots, \hat{\varphi}(a, X_k, t))$ .  $\square$

**Lemma 12.31.** *Suppose that  $a \in \mathbb{K}[X_1, \dots, X_k]$  is well-separating for  $\mathcal{P}$  and such that  $o(\chi(a, T)) = o(\varphi(a, T))$ . Then, for every root  $\tau$  of  $\chi(a, T)$  in  $\mathbb{C}\langle\varepsilon\rangle_b$  such that  $\lim_\varepsilon(\tau) = t$  is a root of multiplicity  $\mu$  of  $\hat{\chi}(a, T)$ ,*

$$\lim_\varepsilon \left( \frac{\varphi(a, X_i, \tau)}{\varphi(a, 1, \tau)} \right) = \frac{\hat{\varphi}^{(\mu-1)}(a, X_i, t)}{\hat{\varphi}^{(\mu-1)}(a, 1, t)}.$$

**Proof:** Let  $\tau_1, \dots, \tau_\ell$  be the roots of  $\chi(a, T)$  in  $\mathbb{C}\langle\varepsilon\rangle_b$  and let  $\tau_{\ell+1}, \dots, \tau_N$  be the roots of  $\chi(a, T)$  in  $\mathbb{C}\langle\varepsilon\rangle \setminus \mathbb{C}\langle\varepsilon\rangle_b$ . Let  $t_j = \lim_\varepsilon(\tau_j)$  for  $j = 1, \dots, \ell$ . Suppose that  $t = \lim_\varepsilon(\tau_1)$  is a root of multiplicity  $\mu$  of  $\hat{\chi}(a, T)$ . By Lemma 12.26 b), there exist  $\mu - 1$  roots of  $\chi(a, T)$  in  $\mathbb{C}\langle\varepsilon\rangle_b$ , numbered  $\tau_2, \dots, \tau_\mu$  without loss of generality, such that

$$t = \lim_\varepsilon(\tau_1) = \dots = \lim_\varepsilon(\tau_\mu),$$

and for every  $\mu < j \leq \ell$ ,  $t_j = \lim_\varepsilon(\tau_j) \neq t$ .

We have

$$\begin{aligned} \chi(a, T) &= \prod_{m \in \{1, \dots, \ell\}} (T - \tau_m) \prod_{m \in \{\ell+1, \dots, N\}} (T - \tau_m), \\ \varphi(a, 1, T) &= \sum_{j=1}^N \prod_{m \in \{1, \dots, N\} \setminus \{j\}} (T - \tau_m), \end{aligned}$$

and  $o(\chi(a, T)) = \sum_{i=\ell+1}^N o(\tau_j)$  by Lemma 12.26. Thus,

$$\begin{aligned} &\varepsilon^{-o(\chi(a, T))} \varphi(a, 1, T) \\ &= \left( \sum_{j=1}^{\mu} \prod_{m \in \{1, \dots, \ell\} \setminus \{j\}} (T - \tau_m) \right) \prod_{m=\ell+1}^N (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \\ &+ \left( \sum_{j=\mu+1}^{\ell} \prod_{m \in \{1, \dots, \ell\} \setminus \{j\}} (T - \tau_m) \right) \prod_{m=\ell+1}^N (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \\ &+ \prod_{m=1}^{\ell} (T - \tau_m) \left( \sum_{j=\ell+1}^N \varepsilon^{-o(\tau_j)} \prod_{m \in J} (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \right). \end{aligned}$$

with  $J = \{\ell + 1, \dots, N\} \setminus \{j\}$ .

Since  $-o(\tau_j) <_\varepsilon 0$ , it follows, taking  $\lim_\varepsilon$ , that

$$\begin{aligned} &\hat{\varphi}(a, 1, T) \\ &= c \mu (T - t)^{\mu-1} \prod_{m=\mu+1}^{\ell} (T - t_m) \\ &+ c \sum_{j=\mu+1}^{\ell} (T - t)^\mu \prod_{m \in \{\mu+1, \dots, N\} \setminus \{j\}} (T - t_m), \end{aligned}$$

with  $c = \prod_{m=\ell+1}^N (-\ln(\tau_j))$ . Thus,

$$\hat{\varphi}^{(\mu-1)}(a, 1, t) = c \mu! \prod_{m=\mu+1}^{\ell} (t - t_m).$$

Denoting by  $\xi_j$  the unique point of  $\text{Zer}(\mathcal{P}, \mathbb{C}\langle\varepsilon\rangle^k)$  such that  $a(\xi_j) = \tau_j$ , and by  $\xi_{ji}$  the  $i$ -th coordinate of  $\xi_j$ , we have similarly

$$\begin{aligned} & \varphi(a, X_i, T) \\ &= \sum_{j=1}^N \xi_{ji} \prod_{m \in \{1, \dots, N\} \setminus \{j\}} (T - \tau_m) \end{aligned}$$

and

$$\begin{aligned} & \varepsilon^{-o(\chi(a, T))} \varphi(a, X_i, T) \\ &= \left( \sum_{j=1}^{\ell} \xi_{ji} \prod_{m \in \{1, \dots, \ell\} \setminus \{j\}} (T - \tau_m) \right) \prod_{m=\ell+1}^N (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \\ &+ \prod_{m=1}^{\ell} (T - \tau_m) C. \end{aligned}$$

with

$$C = \sum_{j=\ell+1}^N \varepsilon^{-o(\tau_j)} \xi_{ji} \prod_{m \in J} (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m).$$

Since  $a$  is well-separating and  $a(\xi_j)$  is bounded, it follows that for all  $i = 1, \dots, \ell, j = 1, \dots, k, \xi_{ji} \in \mathbb{C}\langle\varepsilon\rangle_b$ . So we have  $A \in \mathbb{C}\langle\varepsilon\rangle_b[T]$ , with

$$A = \left( \sum_{j=1}^{\ell} \xi_{ji} \prod_{m \in \{1, \dots, \ell\} \setminus \{j\}} (T - \tau_m) \right) \prod_{m=\ell+1}^N (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m).$$

It is also clear that  $B = \prod_{m=1}^{\ell} (T - \tau_m) \in \mathbb{C}\langle\varepsilon\rangle_b[T]$ .

Since  $\varepsilon^{-o(\chi(a, T))} \varphi(a, X_i, T) \in \mathbb{K}(\varepsilon)_b[T]$ ,  $A \in \mathbb{C}\langle\varepsilon\rangle_b[T]$ ,  $B \in \mathbb{C}\langle\varepsilon\rangle_b[T]$ , and  $B$  is monic,

$$C = \sum_{j=\ell+1}^N \varepsilon^{-o(\tau_j)} \xi_{ji} \prod_{m \in J} (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \in \mathbb{C}\langle\varepsilon\rangle_b[T].$$

So finally

$$\varepsilon^{-o(\chi(a, T))} \varphi(a, X_i, T) = A + \prod_{m=1}^{\mu} (T - \tau_m) \prod_{m=\mu+1}^{\ell} (T - \tau_m) C.$$

Since

$$t = \lim_{\varepsilon} (\tau_1) = \dots = \lim_{\varepsilon} (\tau_{\mu}), t \neq t_j = \lim_{\varepsilon} (\tau_j), j > \mu,$$



and  $a$  is well-separating,

$$\lim_{\varepsilon} (\xi_1) = \dots = \lim_{\varepsilon} (\xi_{\mu}).$$

Denoting by

$$\begin{aligned} x &= (x_1, \dots, x_k) = \lim_{\varepsilon} (\xi_1) = \dots = \lim_{\varepsilon} (\xi_{\mu}), \\ y_j &= (y_{j,1}, \dots, y_{j,k}) = \lim_{\varepsilon} (\xi_j) \end{aligned}$$

and by  $D$  the polynomial obtained by replacing successively  $\varepsilon_m, \varepsilon_{m-1} \dots \varepsilon_1$  by 0 in  $\prod_{m=\mu+1}^{\ell} (T - \tau_m) C$ , we get

$$\begin{aligned} \hat{\varphi}(a, X_i, T) &= c \mu x_i (T - t)^{\mu-1} \prod_{m=\mu+1}^{\ell} (T - t_m) \\ &\quad + (T - t)^{\mu} \left( c \sum_{j=\mu+1}^{\ell} y_{j,i} \prod_{m \in \{\mu+1, \dots, \ell\} \setminus \{j\}} (T - t_m) + D \right), \\ \hat{\varphi}^{(\mu-1)}(a, X_i, t) &= c \mu! x_i \prod_{j=\mu+1}^{\ell} (t - t_j) \end{aligned}$$

with  $c = \prod_{m=\ell+1}^N (-\ln(\tau_j))$ . Finally,

$$\begin{aligned} \lim_{\varepsilon} \left( \frac{\varphi(a, X_i, \tau_j)}{\varphi(a, 1, \tau_j)} \right) &= \lim_{\varepsilon} (\xi_{j,i}) \\ &= x_i \\ &= \frac{\hat{\varphi}^{(\mu-1)}(a, X_i, t)}{\hat{\varphi}^{(\mu-1)}(a, 1, t)}, \end{aligned}$$

for every  $j = 1, \dots, \mu, i = 1, \dots, k$ . □

As is the case for separating elements (see Lemma 4.89), well-separating elements can be chosen in a set defined in advance.

**Lemma 12.32.** *If  $\#\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k) = N$ , then at least one element  $a$  of*

$$\mathcal{A} = \{X_1 + i X_2 + \dots + i^{k-1} X_k \mid 0 \leq i \leq (k-1)N^2\}$$

*is well-separating and such that  $o(\chi(a, T)) = o(\varphi(a, T))$ .*

**Proof:** Define

- $\mathcal{W}_1$ , of cardinality  $\leq N(N-1)/2$ , to be the set of vectors  $\xi - \eta$  with  $\xi$  and  $\eta$  distinct solutions of  $\mathcal{P}$  in  $C\langle\varepsilon\rangle^k$ ,
- $\mathcal{W}_2$ , of cardinality  $\leq N(N-1)/2$ , to be the set of vectors  $\lim_{\varepsilon}(\xi) - \lim_{\varepsilon}(\eta)$  with  $\xi$  and  $\eta$  distinct non-infinitesimally close bounded solutions of  $\mathcal{P}$  in  $C\langle\varepsilon\rangle^k$ ,

- $\mathcal{W}_3$ , of cardinality  $\leq N$ , to be the set of vectors  $c = (c_1, \dots, c_k)$  with  $c_i$  the coefficient of  $\varepsilon^{\max_{i=1, \dots, k} o(\xi_i)}$  in  $\xi_i$ , for  $\xi = (\xi_1, \dots, \xi_k) \in \text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$ ,
- $\mathcal{W} = \mathcal{W}_1 \cup \mathcal{W}_2 \cup \mathcal{W}_3$ . Note that  $\mathcal{W}$  is of cardinality  $\leq N^2$

If  $j$  is such that, for every  $c \in \mathcal{W}_3$ ,  $c_1 + \dots + j^{k-1} c_k \neq 0$  and

$$a = X_1 + \dots + j^{k-1} X_k,$$

then for every  $\xi \in \text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$ ,

$$\text{In}(a(\xi)) = c_1 + \dots + j^{k-1} c_k$$

and  $o(a(\xi)) = \max_{i=1, \dots, k} (o(\xi_i))$ .

If, for every  $\xi = (\xi_1, \dots, \xi_k) \in \text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$ ,  $o(a(\xi)) = \max_{i=1, \dots, k} (o(\xi_i))$  then  $a$  maps unbounded elements of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  to unbounded elements of  $C\langle\varepsilon\rangle$ . Denote by  $\xi_1, \dots, \xi_\ell$  the elements of  $\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)$  and by  $\xi_{\ell+1}, \dots, \xi_N$  the elements of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k) \setminus \text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)$ . Then  $\tau_1 = a(\xi_1), \dots, \tau_\ell = a(\xi_\ell)$  are the roots of  $\chi(a, T)$  in  $C\langle\varepsilon\rangle_b$  and  $\tau_{\ell+1} = a(\xi_{\ell+1}), \dots, \tau_N = a(\xi_N)$  are the roots of  $\chi(a, T)$  in  $C\langle\varepsilon\rangle \setminus C\langle\varepsilon\rangle_b$ . For  $i = 1, \dots, k$ ,

$$\begin{aligned} & \varepsilon^{-o(\chi(a, T))} \varphi(a, X_i, T) \\ &= \left( \sum_{j=1}^{\ell} \xi_{ji} \prod_{m \in \{1, \dots, \ell\} \setminus \{j\}} (T - \tau_m) \right) \prod_{m=\ell+1}^N (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m) \\ &+ \prod_{m=1}^{\ell} (T - \tau_m) \left( \sum_{j=\ell+1}^N \varepsilon^{-o(\tau_j)} \xi_{ji} C \right) \end{aligned}$$

with

$$C = \sum_{j=\ell+1}^N \varepsilon^{-o(\tau_j)} \xi_{ji} \prod_{m \in \{\ell+1, \dots, N\} \setminus \{j\}} (\varepsilon^{-o(\tau_m)} T - \varepsilon^{-o(\tau_m)} \tau_m)$$

belongs to  $K(\varepsilon)_b[T]$ , since  $o(\tau_j) = o(a(\xi_j)) \geq_\varepsilon o(\xi_{j,i})$  and  $o(\varepsilon^{-o(\tau_j)} \xi_{ji}) \leq_\varepsilon 0$ .

So, if  $j$  is such that, for every  $w \in \mathcal{W}$ ,  $w_1 + \dots + j^{k-1} w_k \neq 0$ , then  $a = X_1 + \dots + j^{k-1} X_k$  is well-separating and such that

$$o(\chi(a, T)) = o(\phi(a, T)).$$

For a fixed  $w \in \mathcal{W}$ , there are at most  $k - 1$  elements of  $\mathcal{A}$  that satisfy  $w_1 + \dots + j^{k-1} w_k = 0$ . This is because an element

$$X_1 + j X_2 + \dots + j^{k-1} X_k$$

satisfying  $w_1 + \dots + j^{k-1} w_k = 0$  is such that  $P_w(j) = 0$ , with

$$P_w(T) = w_1 + T w_2 + \dots + T^{k-1} w_k.$$

But  $P_w(T)$ , which is non-zero, has at most  $k - 1$  roots. So the result is clear by the pigeon-hole principle.  $\square$

According to the preceding results, the set  $\lim_\varepsilon (\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k) \cap \mathbb{R}^k$  can be obtained as follows:

- Determine a well-separating element  $a = X_1 + \dots + j^{k-1}X_k$  such that  $o(\chi(a, T)) = o(\varphi(a, T))$  as follows:
  - List all  $a \in \mathcal{A}$  that are separating and compute the corresponding  $\hat{\chi}(a, T)$ .
  - Among these list the  $a$  such that the degree of  $\hat{\chi}(a, T)$  is minimal. This condition guarantees that  $a$  maps unbounded elements of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  to unbounded roots of  $\chi(a, T)$  since, by Corollary 12.27,  $\deg(\chi(a, T)) - \deg(\hat{\chi}(a, T))$  is the number of unbounded roots of  $\chi(a, T)$  and is maximal when all unbounded elements of  $\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k)$  have unbounded images by  $a$ .
  - Among these list those such that  $o(\chi(a, T)) = o(\varphi(a, T))$ .
  - Among these find an  $a$  such that the number of distinct roots of  $\hat{\chi}(a, T)$  is maximal, i.e. such that  $\deg(\gcd(\hat{\chi}(a, T), \hat{\chi}'(a, T)))$  is minimal. This guarantees that no two non-infinitesimally close elements of  $\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k)$  are sent by  $a$  to two infinitesimally close elements of  $C\langle\varepsilon\rangle_b$ .
- Lemma 12.32 guarantees that there exists such an  $a$  in  $\mathcal{A}$ .
- For such an  $a$  and every root  $t$  of  $\hat{\chi}(a, T)$  in  $\mathbb{R}$  with multiplicity  $\mu$  consider

$$x_i = \frac{\hat{\varphi}^{(\mu-1)}(a, X_i, t)}{\hat{\varphi}^{(\mu-1)}(a, 1, t)}.$$

The root of  $\hat{\chi}(a, T)$  in  $\mathbb{R}$  can be described by its Thom encoding. All elements of  $\lim_\varepsilon (\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k) \cap \mathbb{R}^k$  are obtained this way since if  $x \in \lim_\varepsilon (\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k) \cap \mathbb{R}^k$ ,  $a(x) \in \mathbb{R}$  is a root of  $\hat{\chi}(a, T)$ . Conversely if  $t \in \mathbb{R}$  is a root of  $\hat{\chi}(a, T)$  of multiplicity  $\mu$ ,

$$x_i = \frac{\hat{\varphi}^{(\mu-1)}(a, X_i, t)}{\hat{\varphi}^{(\mu-1)}(a, 1, t)} \in \mathbb{R},$$

since  $\hat{\varphi}(a, X_i, T) \in K[T]$ ,  $t \in \mathbb{R}$ .

We can now describe an algorithm for computing the limit of the bounded solutions of a polynomial system. Since we want to perform the computations in a ring rather than in a field, the following remark will be useful.

*Remark 12.33.* Let  $\#\text{Zer}(\mathcal{P}, C\langle\varepsilon\rangle^k) = N$ . Consider  $b \neq 0 \in K(\varepsilon)$ . Using Notation 12.30, we have

$$(\chi(ba, bT), \varphi_b(ba, bT)) = b^N (\chi(a, T), \varphi(a, T)).$$

Thus,  $o(\chi(a, T)) = o(\varphi(a, T))$  if and only if  $o(\chi(b a, b T)) = o(\varphi_b(b a, b T))$ .  $\square$

*Algorithm 12.14. [Limits of Bounded Points]*

- **Structure:** an ordered ring  $D$  with division in  $\mathbb{Z}$  contained in an ordered field  $K$ .
- **Input:**  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$ , a zero-dimensional polynomial system with only simple zeroes  $\mathcal{P} \subset D[\varepsilon][X_1, \dots, X_k]$ , a basis  $\mathcal{B}$  of

$$A = K(\varepsilon)[X_1, \dots, X_k] / \text{Ideal}(\mathcal{P}, K(\varepsilon))$$

such that the multiplication table  $\mathcal{M}$  of  $A$  in  $\mathcal{B}$  has entries in  $D[\varepsilon]$ , an element  $b \neq 0 \in D[\varepsilon]$  such that  $b, b X_1, \dots, b X_k$  have coordinates in  $D$  in the basis  $\mathcal{B}$ .

- **Output:** a set  $\mathcal{U}$  of real univariate representations, such that the set of points in  $\mathbb{R}^k$  associated to these  $k + 2$  tuples are the elements of  $\lim_{\varepsilon} (\text{Zer}_b(\mathcal{P}, C\langle\varepsilon\rangle^k) \cap \mathbb{R}^k)$ .
- **Complexity:**

$$(\lambda d_1 \dots d_k)^{O(m)}$$

when  $\mathcal{P}$  is a special Gröbner basis

$$\mathcal{P} = \{b_1 X_1^{d_1} + Q_1, \dots, b_k X_k^{d_k} + Q_k\} \subset D[\varepsilon][X_1, \dots, X_k],$$

$$\deg_X(Q_i) < d_i, \deg_{X_j}(Q_i) < d_j \quad \deg_{\varepsilon}(Q_i) \leq \lambda, \deg_{\varepsilon}(b_i) \leq \lambda, b = b_1 \dots b_k.$$

- **Procedure:**
  - For every  $a = X_1 + j X_2 + \dots + j^{k-1} X_k, j = 0, \dots, (k - 1) N^2$  compute  $(\chi(b a, T), \varphi_b(b a, b, T))$  using Algorithm 12.11 (Candidate Univariate Representation).
  - Keep the values of  $a$  such that  $o(\chi(b a, b T)) = o(\varphi_b(b a, b T))$ .
  - Compute

$$\begin{aligned} \hat{\chi}(b a, b T) &= \lim_{\varepsilon} (\varepsilon^{-o(\chi(b a, b T))} \chi(b a, b T)) \\ \hat{\varphi}(b a, b T) &= \lim_{\varepsilon} (\varepsilon^{-o(\chi(b a, b T))} c \varphi_b(b a, b T)). \end{aligned}$$

- Choose an  $a$  among those for which  $\deg(\hat{\chi}(b a, b T))$  is minimal and for which  $\deg(\text{gcd}(\hat{\chi}(b a, b T), \hat{\chi}'(b a, b T)))$  is minimal (Notation 12.30), computing  $\text{gcd}(\hat{\chi}(b a, b T), \hat{\chi}'(b a, b T))$  using Remark 10.18.
- Return  $(\hat{\chi}(b a, b T), \hat{\varphi}(b a, b T))$ .
- Compute the list of Thom encodings of the roots of  $\hat{\chi}(b a, b T)$  in  $\mathbb{R}$  using Algorithm 10.14 (Thom Encoding) and Remark 10.76. Read from the Thom encoding  $\sigma$  the multiplicity  $\mu$  of the associated root  $t_{\sigma}$ . For every such Thom encoding  $\sigma$ , place  $(\hat{\chi}(b a, b T), \hat{\varphi}^{(\mu-1)}(b a, b T), \sigma)$  in  $\mathcal{U}$ .

**Proof of correctness:** The correctness follows from the discussion preceding the algorithm, using Remark 12.33.  $\square$

**Complexity analysis:** We estimate the complexity only in the case where  $\mathcal{P}$  is a special Gröbner basis contained in  $D[\varepsilon][X_1, \dots, X_k]$ , since this is the only way we are going to use it later. Let

$$\mathcal{P} = \{b_1 X_1^{d_1} + Q_1, \dots, b_k X_k^{d_k} + Q_k\} \subset D[\varepsilon][X_1, \dots, X_k],$$

with  $\deg_X(Q_i) < d_i$ ,  $\deg_{X_j}(Q_i) < d_j$ ,  $\deg_\varepsilon(Q_i) \leq \lambda$ ,  $\deg_\varepsilon(b_i) \leq \lambda$ ,  $b = b_1 \cdots b_k$ .

Then the number of arithmetic operations in  $D[\varepsilon]$  is  $(d_1 \cdots d_k)^{O(1)}$  according to the complexity analysis of Algorithm 12.11 (Univariate Representation), Remark 10.18, and Remark 10.76. The degrees in  $\varepsilon$  of the polynomials occurring in the multiplication table are

$$\lambda k d_1 (d_1 + \cdots + d_{k-1} - k + 1)$$

according to the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table). Finally, using the complexity of Algorithm 8.4 (Addition of multivariate polynomials) and Algorithm 8.5 (Multiplication of multivariate polynomials), the complexity in  $D$  is

$$(\lambda d_1 \cdots d_k)^{O(m)}.$$

The degree in  $T$  and number of real univariate representations output is  $d_1 \cdots d_k$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring in the computation of the multiplication table and its output are

$$(d_1 + \cdots + d_{k-1} - k + 1) (k d_1 + 1) (\tau + 4 \nu'),$$

where  $\nu'$  is the bitsize of  $(\lambda (d_1 + \cdots + d_k) + 1)^m$ , according to the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).  $\square$

In later chapters of the book, we need a parametrized version of this algorithm in the case of a parametrized special system.

A **parametrized univariate representation** with parameters  $Y$  is a  $k + 2$ -tuple

$$\begin{aligned} u(Y) &= (f(Y, T), g(Y, T)) \in D[Y][T]^{k+2}, \\ g(Y, T) &= (g_0(Y, T), g_1(Y, T), \dots, g_k(Y, T)). \end{aligned}$$

We need a notation. Let  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$ . If  $f \in A[\varepsilon]$  and  $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{N}^m$ , we denote by  $f_\alpha \in A$  the coefficient of  $\varepsilon^\alpha$  in  $f$  and by  $g_\alpha = (g_{0\alpha}, g_{1\alpha}, \dots, g_{k\alpha})$

*Algorithm 12.15. [Parametrized Limit of Bounded Points]*

- **Structure:** an ordered ring  $D$  with division in  $\mathbb{Z}$  contained in an ordered field  $K$ .
- **Input:**  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$ , a parametrized special Gröbner basis

$$\mathcal{P} = \{b_1 X_1^{d_1} + Q_1(Y, X), \dots, b_k X_k^{d_k} + Q_k(Y, X)\} \subset D[Y, \varepsilon][X_1, \dots, X_k]$$

with  $Y = (Y_1, \dots, Y_\ell)$ , and the corresponding parametrized multiplication table  $\mathcal{M}$  with entries in  $D[Y, \varepsilon]$ , with  $b_i \in D[\varepsilon]$ .

- **Output:** a set  $\mathcal{U}$  of parametrized univariate representations of the form,  $u(Y) = (f, g)$ , where  $(f, g) \in D[Y][T]^{k+2}$ . The set  $\mathcal{U}$  has the property that for any point  $y \in \mathbb{R}^\ell$ , denoting by  $\mathcal{U}(y)$  the subset of  $\mathcal{U}$  such that  $f(y, T)$  and  $g_0(y, T)$  are coprime, the points associated to the univariate representations  $u(y)$  in  $\mathcal{U}(y)$  contain  $\lim_\varepsilon (\text{Zer}_b(\mathcal{G}_y, C\langle \varepsilon \rangle^k)) \cap \mathbb{R}^k$ .
- **Complexity:**  $((\lambda + t) d_1 \dots d_k)^{O(m+\ell)}$ .
- **Procedure:**
  - $b := b_1 \dots b_k$ .
  - For every

$$a = X_1 + j X_2 + \dots + j^{k-1} X_k, \quad j = 0, \dots, (k - 1)N^2,$$

compute the parametrized univariate representation

$$(\chi(ba, bT), \varphi_b(ba, bT))$$

by performing the computations of Algorithm 12.11 (Candidate Univariate Representation) in  $D[Y, \varepsilon]$ .

- For every  $\alpha \in \mathbb{Z}^m$  such that  $\varepsilon^\alpha$  appears in  $\chi$ , for every  $\mu \leq \deg_T(\chi_\alpha)$ , include  $(\chi_\alpha(ba, bT), \varphi_{b, \alpha}^{(\mu-1)}(ba, bT))$  in the set  $\mathcal{U}$ .
- Output  $\mathcal{U}$ .

**Proof of correctness:** The correctness follows from the discussion preceding Algorithm 12.14. In this parametric situation, the choice of the well-separating element and the order of the univariate representation depends on the parameters, as well as the multiplicities of the roots. This is the reason why we place all the possibilities in  $\mathcal{U}$ . □

**Complexity analysis:** Let

$$\mathcal{P} = \{b_1 X_1^{d_1} + Q_1, \dots, b_k X_k^{d_k} + Q_k\} \subset D[Y][\varepsilon][X_1, \dots, X_k],$$

with  $\deg_X(Q_i) < d_i$ ,  $\deg_{X_j}(Q_i) < d_j$ ,  $\deg_\varepsilon(Q_i) \leq \lambda$ ,  $\deg_Y(Q_i) \leq t$ , and  $\deg_\varepsilon(b_i) \leq \lambda$ ,  $b = b_1 \dots b_k$ . The number of arithmetic operations in  $D[Y][\varepsilon]$  is  $(d_1 \dots d_k)^{O(1)}$  according to the complexity analysis of Algorithm 12.11 (Univariate Representation), Remark 10.18 and Remark 10.76. The degrees in  $\varepsilon$  and  $Y$  of the polynomials occurring in the multiplication table are respectively  $O(\lambda k d_1 (d_1 + \dots + d_k))$  and  $O(t k d_1 (d_1 + \dots + d_k))$  according to the complexity analysis of 12.10 (Parametrized Special Multiplication Table). Finally, using the complexity of Algorithm 8.4 (Addition of multivariate polynomials) and Algorithm 8.5 (Multiplication of multivariate polynomials), the complexity in  $D$  is

$$((\lambda + t) d_1 \dots d_k)^{O(m+\ell)}.$$

The degrees in  $T$  of the of real univariate representations output is  $d_1 \cdots d_k$ , and their degrees in  $Y$  is  $(d_1 \cdots d_k)^{O(1)}$ . The number of real univariate representations output is  $(d_1 \cdots d_k)^{O(1)}$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring computing the multiplication table and its output are  $(d_1 + \cdots + d_{k-1} - k + 1)(k d_1 + 1)(\tau + 4\nu')$ , where  $\nu'$  is the bitsize of  $((\lambda + t)(d_1 + \cdots + d_k) + 1)^{m+\ell}$ , according to the complexity analysis of 12.10 (Parametrized Special Multiplication Table).  $\square$

## 12.6 Finding Points in Connected Components of Algebraic Sets

We are going to describe a method for finding at least one point in every semi-algebraically connected component of an algebraic set. We know by Proposition 7.9 that when we consider a bounded nonsingular algebraic hypersurface, it is possible to change coordinates so that its projection to the  $X_1$ -axis has a finite number of non-degenerate critical points. These points provide at least one point in every semi-algebraically connected component of the bounded nonsingular algebraic hypersurface by Proposition 7.4. Unfortunately this result is not very useful in algorithms since it provides no method for performing this linear change of variables. Moreover when we deal with the case of a general algebraic set, which may be unbounded or singular, this method no longer works.

We first explain how to associate to a possibly unbounded algebraic set  $Z \subset \mathbb{R}^k$  a bounded algebraic set  $Z' \subset \mathbb{R}\langle \varepsilon \rangle^{k+1}$ , whose semi-algebraically connected components are closely related to those of  $Z$ .

Let  $Z = \text{Zer}(Q, \mathbb{R}^k)$  and consider

$$Z' = \text{Zer}(Q^2 + (\varepsilon^2(X_1^2 + \cdots + X_{k+1}^2) - 1)^2, \mathbb{R}\langle \varepsilon \rangle^{k+1}).$$

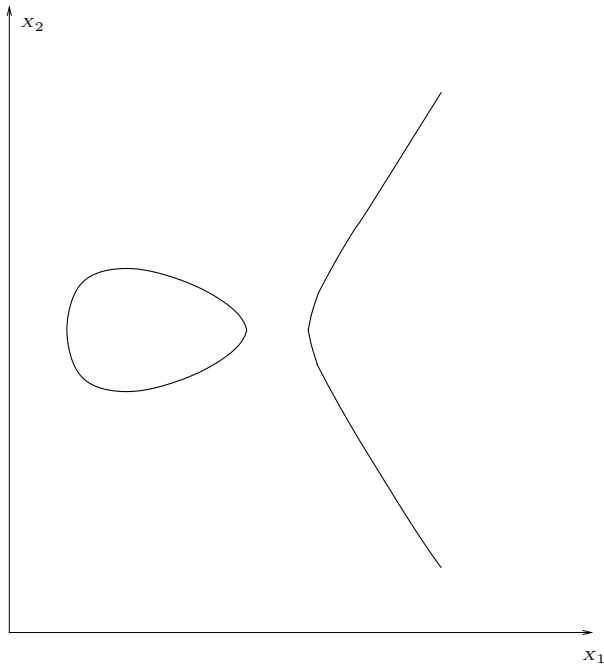
The set  $Z'$  is the intersection of the sphere  $S_\varepsilon^k$  of center 0 and radius  $1/\varepsilon$  with a cylinder based on the extension of  $Z$  to  $\mathbb{R}\langle \varepsilon \rangle$ . The intersection of  $Z'$  with the hyperplane  $X_{k+1} = 0$  is the intersection of  $Z$  with the sphere  $S_\varepsilon^{k-1}$  of center 0 and radius  $1/\varepsilon$ . Denote by  $\pi$  the projection from  $\mathbb{R}\langle \varepsilon \rangle^{k+1}$  to  $\mathbb{R}\langle \varepsilon \rangle^k$ .

**Proposition 12.34.** *Let  $N$  be a finite set of points meeting every semi-algebraically connected component of  $Z'$ . Then  $\pi(N)$  meets every semi-algebraically connected component of the extension  $\text{Ext}(Z, \mathbb{R}\langle \varepsilon \rangle)$  of  $Z$  to  $\mathbb{R}\langle \varepsilon \rangle$ .*

**Proof:** Let  $D$  a semi-algebraically connected components of  $Z$ . If  $D$  is bounded,  $\text{Ext}(D, \mathbb{R}\langle \varepsilon \rangle)$  does not intersect  $S_\varepsilon^{k-1}$ , and  $\pi^{-1}(\text{Ext}(D, \mathbb{R}\langle \varepsilon \rangle))$  is semi-algebraically homeomorphic to two copies of  $\text{Ext}(D, \mathbb{R}\langle \varepsilon \rangle)$ , one in each of the half-spaces defined, respectively, by  $X_{k+1} > 0$  and by  $X_{k+1} < 0$ . Thus, since  $N$  intersects every semi-algebraically connected component of  $Z'$ ,  $N$  intersects  $\pi^{-1}(\text{Ext}(D, \mathbb{R}\langle \varepsilon \rangle))$  and  $\pi(N)$  intersects  $\text{Ext}(D, \mathbb{R}\langle \varepsilon \rangle)$ .

If  $D$  is unbounded, the set  $A$  of elements  $r \in \mathbb{R}$  such that  $D$  intersects the sphere  $S^{k-1}(0, r)$  of center 0 and radius  $r$  is semi-algebraic and unbounded and contains an open interval  $(a, +\infty)$ . Thus  $1/\varepsilon \in \text{Ext}(A, \mathbb{R}\langle\varepsilon\rangle)$ , and  $\text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$  intersects  $S_\varepsilon^{k-1}$ . Take  $z \in \text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle) \cap S_\varepsilon^{k-1}$ , and denote by  $D'$  the semi-algebraically connected component of  $Z'$  containing  $z' = (z, 0) \in Z'$ . Take  $x \in D' \cap N$  and consider a semi-algebraic path  $\gamma$  connecting  $z'$  to  $x$  inside  $D'$ . Then,  $\pi(\gamma)$  is a semi-algebraic path connecting  $z$  to  $\pi(x)$  inside  $\text{Ext}(Z, \mathbb{R}\langle\varepsilon\rangle)$ , thus  $\pi(x)$  and  $z$  belong to the same semi-algebraically connected component of  $\text{Ext}(Z, \mathbb{R}\langle\varepsilon\rangle)$ . Since  $z \in \text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$ , then  $\pi(x) \in \text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$ , and  $\pi(N)$  intersects  $\text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$ .  $\square$

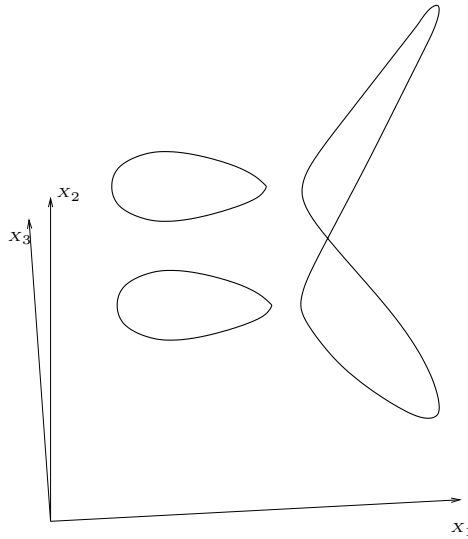
Let us illustrate this result. If  $Q = X_2^2 - X_1(X_1 - 1)(X_1 + 1)$ , then  $Z = \text{Zer}(Q, \mathbb{R}^2)$  is a cubic curve with one bounded semi-algebraically connected component and one unbounded semi-algebraically connected component (see Figure 12.3).



**Fig. 12.3.** Cubic curve in the plane

The corresponding  $Z' \subset \mathbb{R}\langle\varepsilon\rangle^3$  (see Figure 12.4) has two semi-algebraically connected components above the bounded semi-algebraically connected component of the cubic curve, and one semi-algebraically connected component above the unbounded semi-algebraically connected component of the cubic curve.





**Fig. 12.4.** Cubic curve lifted to a big sphere

So, if we have a method for finding a point in every semi-algebraically connected component of a bounded algebraic set, we obtain immediately, using Proposition 12.34, a method for finding a point in every connected component of an algebraic set. Note that these points have coordinates in the extension  $\mathbb{R}(\varepsilon)$  rather than in the real closed field  $\mathbb{R}$  we started with. However, the extension from  $\mathbb{R}$  to  $\mathbb{R}(\varepsilon)$  preserves semi-algebraically connected components (Proposition 5.24).

We are going to define  $X_1$ -pseudo-critical points of  $\text{Zer}(Q, \mathbb{R}^k)$  when  $\text{Zer}(Q, \mathbb{R}^k)$  is a bounded algebraic set. These pseudo-critical points are a finite set of points meeting every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ . They are the limits of the critical points of the projection to the  $X_1$  coordinate of a bounded nonsingular algebraic hypersurface defined by a particular infinitesimal perturbation of the polynomial  $Q$ . Moreover, the equations defining the critical points of the projection on the  $X_1$  coordinate on the perturbed algebraic set have the special algebraic structure considered in Proposition 12.7.

Given a polynomial  $Q \in \mathbb{R}[X_1, \dots, X_k]$  we define  $\text{tDeg}_{X_i}(Q)$ , the **total degree of  $Q$  in  $X_i$** , as the maximal total degree of the monomials in  $Q$  containing the variable  $X_i$ .

**Notation 12.35. [Deformation]** Let  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ , and  $c$

$$\begin{aligned} G_k(\bar{d}, c) &= c^{\bar{d}_1} (X_1^{\bar{d}_1} + \dots + X_k^{\bar{d}_k} + X_2^2 + \dots + X_k^2) - (2k - 1), \\ \text{Def}(Q, \zeta) &= \zeta G_k(\bar{d}, c) + (1 - \zeta) Q. \end{aligned} \tag{12.17}$$

□

In the next pages, the polynomial  $Q \in D[X_1, \dots, X_k]$ , where  $D$  is a ring contained in the real closed field  $R$ , and  $(d_1, \dots, d_k)$  satisfy the following conditions:

- $Q(x) \geq 0$  for every  $x \in R^k$ ,
- $\text{Zer}(Q, R^k) \subset B(0, 1/c)$  for some  $c \leq 1, c \in D$ ,
- $d_1 \geq d_2 \geq \dots \geq d_k$ ,
- $\text{deg}(Q) \leq d_1, \text{tDeg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ .

Note that  $\forall x \in B(0, 1/c) \quad G_k(\bar{d}, c)(x) < 0$ .

*Remark 12.36.* Note that supposing  $Q(x) \geq 0$  for every  $x \in R^k$  is not a big loss of generality since we can always replace  $Q$  by  $Q^2$  if it is not the case. Note also that we can always take

$$d_1 = \dots = d_k = \text{deg}(Q).$$

However considering different  $d_i$  will be useful when the degree with respect to some variables is small. □

Let  $\bar{d}_i$  be an even number  $> d_i, i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ .

Let  $\zeta$  be a variable and  $R\langle\zeta\rangle$  be as usual the field of algebraic Puiseux series in  $\zeta$  with coefficients in  $R$ .

**Proposition 12.37.**

$$\lim_{\zeta} (\text{Zer}(\text{Def}(Q, \zeta), R\langle\zeta\rangle^k)) = \text{Zer}(Q, R^k).$$

Moreover  $\text{Zer}(\text{Def}(Q, \zeta), R\langle\zeta\rangle^k) \subset B(0, 1/c)$ .

**Proof:** Since  $\lim_{\zeta}$  is a ring homomorphism from  $R\langle\zeta\rangle_b$  to  $R$ , it is clear that  $\lim_{\zeta}(\text{Zer}(\text{Def}(Q, \zeta), R\langle\zeta\rangle^k)) \subset \text{Zer}(Q, R^k)$ . We show that

$$\text{Zer}(Q, R^k) \subset \lim_{\zeta}(\text{Zer}(\text{Def}(Q, \zeta), R\langle\zeta\rangle^k)).$$

Let  $x \in \text{Zer}(Q, R^k)$ . Since  $\text{Zer}(Q, R^k)$  is bounded, for every  $r > 0$  in  $R$  there is a  $y \in B(x, r)$  such that  $Q(y) > 0$ . Thus, using Theorem 3.19 (Curve selection lemma), there exists a semi-algebraic path  $\gamma$  from  $[0, 1]$  to  $R^k$  starting from  $x$  such that  $Q(\gamma(t)) \neq 0$  for  $t \in (0, 1]$ . By Theorem 3.20, the set  $\gamma([0, 1])$  is a bounded subset of  $R^k$ . Denote by  $\bar{\gamma}$  the extension of  $\gamma$  to  $R\langle\zeta\rangle$ , and note that  $\bar{\gamma}([0, 1])$  is a bounded subset of  $R\langle\zeta\rangle^k$ , using Proposition 2.87. Since  $\text{Def}(Q, \zeta)(\bar{\gamma}(0)) < 0$  and  $\text{Def}(Q, \zeta)(\bar{\gamma}(t)) > 0$ , for every  $t \in R$ , with  $0 < t < 1$ , there exists  $\tau \in R\langle\zeta\rangle, \lim_{\zeta}(\tau) = 0$ , such that  $\text{Def}(Q, \zeta)(\bar{\gamma}(\tau)) = 0$  by Proposition 3.4. Since  $\lim_{\zeta}$  is a ring homomorphism

$$Q(\lim_{\zeta}(\bar{\gamma}(\tau))) = \lim_{\zeta}(Q(\bar{\gamma}(\tau))) = \lim_{\zeta}(\text{Def}(Q, \zeta)(\bar{\gamma}(\tau))) = 0.$$

Since  $\lim_{\zeta}(\tau) = 0, \gamma(0) = x$ , and  $\gamma$  is continuous, we have

$$\lim_{\zeta}(\bar{\gamma}(\tau)) = \gamma(0) = x,$$

using Lemma 3.21. Thus we have found

$$\bar{\gamma}(\tau) \in \text{Zer}(\text{Def}(Q, c), \mathbb{R}\langle \zeta \rangle^k)$$

such that  $\lim_{\zeta} (\bar{\gamma}(\tau)) = x$ .

Since

$$\lim_{\zeta} (\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)) = \text{Zer}(Q, \mathbb{R}^k)$$

and every point  $x = (x_1, \dots, x_k) \in \text{Zer}(Q, \mathbb{R}^k)$  satisfies  $x_1^2 + \dots + x_k^2 < 1/c$ , it follows clearly that every point  $y = (y_1, \dots, y_k) \in \text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  satisfies  $y_1^2 + \dots + y_k^2 < 1/c$ .  $\square$

**Proposition 12.38.** *The algebraic set  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  is a non-singular algebraic hypersurface bounded over  $\mathbb{R}$ .*

**Proof:** The fact that  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  is bounded follows from Proposition 12.37.

To prove that  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  is a non-singular hypersurface, consider the function

$$\Phi(x) = \frac{Q(x)}{Q(x) - G_k(\bar{d}, c)(x)}$$

from  $\mathbb{R}^k \setminus \text{Zer}(Q - G_k(\bar{d}, c), \mathbb{R}^k)$  to  $\mathbb{R}$ . By Sard's Theorem (Theorem 5.56) the set of critical values of  $\Phi$  is finite. So there is an  $a \in \mathbb{R}$ ,  $a > 0$ , such that for every  $b \in (0, a)$  the function  $\Phi$  has no critical value.

Since  $\text{Zer}(\text{Def}(Q, b), \mathbb{R}^k) \cap \text{Zer}(Q - G_k(\bar{d}, c), \mathbb{R}^k) = \emptyset$ ,

$$\text{Zer}(\text{Def}(Q, b), \mathbb{R}^k) = \{x \in \mathbb{R}^k \mid \Phi(x) = b\}.$$

The set  $\text{Zer}(\text{Def}(Q, b), \mathbb{R}^k)$  is a non-singular algebraic hypersurface, since  $\text{Grad}(\text{Def}(Q, b))(x) = 0$  on  $\text{Zer}(\text{Def}(Q, b), \mathbb{C}^k)$  implies that  $\text{Grad}(\Phi)(x) = 0$ . So the formula  $\Psi(a)$  defined by

$$\forall b \forall x (0 < b < a \wedge \text{Def}(Q, b)(x) = 0) \Rightarrow \text{Grad}(\text{Def}(Q, b))(x) \neq 0$$

is true in  $\mathbb{R}$ . Using Theorem 2.80 (Tarski-Seidenberg principle),  $\Psi(a)$  is true in  $\mathbb{R}\langle \zeta \rangle$  which contains  $\mathbb{R}$ . Hence, since  $0 < \zeta < a$ ,  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  is a non-singular algebraic hypersurface.  $\square$

**Notation 12.39.** Let  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ , and using Notation 12.35, consider

$$\begin{aligned} \text{Cr}(Q, \zeta) &= \left\{ \text{Def}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k} \right\}, \\ \text{Def}_+(Q, \zeta) &= \text{Def}(Q, \zeta) + X_{k+1}^2, \\ \text{Cr}_+(Q, \zeta) &= \left\{ \text{Def}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k}, 2X_{k+1} \right\}. \end{aligned}$$

$\square$

Note that

$$\text{Zer}(\text{Cr}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$$

is the set of  $X_1$ -critical points on

$$\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$$

i.e. the critical points on  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  of the projection map to the  $X_1$  coordinate.

The following lemma is easy to prove using the arguments in the proofs of Propositions 12.38, and 12.44.

**Lemma 12.40.** *The algebraic set  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  is a non-singular algebraic hypersurface which is bounded over  $\mathbb{R}$ . Moreover,*

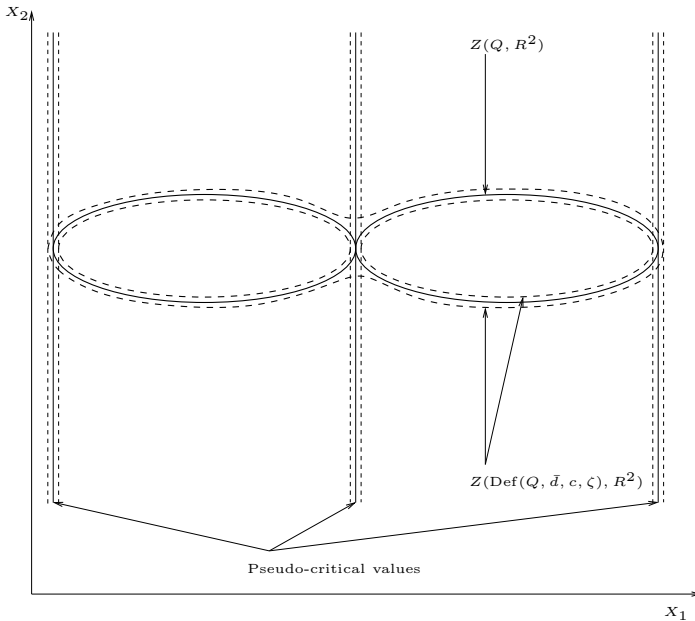
$$\lim_{\zeta} (\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})) = \text{Zer}(Q, \mathbb{R}^k) \times \{0\},$$

and  $\pi$  (the projection of  $(x_1, \dots, x_{k+1}) \in \mathbb{R}\langle \zeta \rangle^{k+1}$  to  $x_1 \in \mathbb{R}\langle \zeta \rangle$ ) has a finite number of critical points on  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$ .

Note that an  $X_1$ -critical point on  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  must have its last coordinate 0 and thus its first  $k$  coordinates define an  $X_1$ -critical point on  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$ .

**Definition 12.41.** An  $X_1$ -pseudo-critical point on  $\text{Zer}(Q, \mathbb{R}^k)$  is the  $\lim_{\zeta}$  of an  $X_1$ -critical point on  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$ .

An  $X_1$ -pseudo-critical value on  $\text{Zer}(Q, \mathbb{R}^k)$  is the projection to the  $X_1$ -axis of an  $X_1$ -pseudo-critical point on  $\text{Zer}(Q, \mathbb{R}^k)$ . □



**Fig. 12.5.** Pseudo-critical values of an algebraic set in  $\mathbb{R}^2$

According to Definition 12.41, an  $X_1$ -pseudo-critical point of  $\text{Zer}(Q, \mathbb{R}^k)$  is the  $\lim_\zeta$  of an  $X_1$ -critical point on  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle\zeta\rangle^k)$ , so that an  $X_1$ -pseudo-critical point on  $\text{Zer}(Q, \mathbb{R}^k)$  is also the  $\lim_\zeta$  of an  $X_1$ -critical point on  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle\zeta\rangle^{k+1})$ .

**Proposition 12.42.** *The set of  $X_1$ -pseudo-critical points on  $\text{Zer}(Q, \mathbb{R}^k)$  meets every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ .*

The proof of Proposition 12.42 will use the following result.

**Proposition 12.43.** *If  $S' \subset \mathbb{R}\langle\zeta\rangle^k$  is a semi-algebraic set, then  $\lim_\zeta(S')$  is a closed semi-algebraic set. Moreover, if  $S' \subset \mathbb{R}\langle\zeta\rangle^k$  is a semi-algebraic set bounded over  $\mathbb{R}$  and semi-algebraically connected, then  $\lim_\zeta(S')$  is semi-algebraically connected.*

**Proof:** Using Proposition 2.82, we can suppose that  $S' \subset \mathbb{R}\langle\zeta\rangle^k$  is described by a quantifier free formula  $\Phi(X, \zeta)$  with coefficients in  $\mathbb{R}[\zeta]$ . Introduce a new variable  $X_{k+1}$  and denote by  $\Phi(X, X_{k+1})$  the result of substituting  $X_{k+1}$  for  $\zeta$  in  $\Phi(X, \zeta)$ . Embed  $\mathbb{R}^k$  in  $\mathbb{R}^{k+1}$  by sending  $X$  to  $(X, 0)$ .

We prove that  $\lim_\zeta(S') = \overline{T} \cap \text{Zer}(X_{k+1}, \mathbb{R}^{k+1})$ , where

$$T = \{(x, x_{k+1}) \in \mathbb{R}^{k+1} \mid \Phi(x, x_{k+1}) \wedge x_{k+1} > 0\}$$

and  $\overline{T}$  is the closure of  $T$ . If  $x \in \lim_\zeta(S')$ , then there exists  $z \in S'$  such that  $\lim_\zeta(z) = x$ . Since  $(z, \zeta)$  belongs to the extension of  $B(x, r) \cap T$  to  $\mathbb{R}\langle\zeta\rangle$ , it follows that  $B(x, r) \cap T$  is non-empty for every  $r \in \mathbb{R}$ ,  $r > 0$ , and hence that  $x \in \overline{T}$ . Conversely, let  $x$  be in  $\overline{T} \cap \text{Zer}(X_{k+1}, \mathbb{R}^{k+1})$ . For every  $r \in \mathbb{R}$ , with  $r > 0$ ,  $B(x, r) \cap T \cap \text{Zer}(X_{k+1}, \mathbb{R}^{k+1})$  is non-empty, and hence, according to Theorem 2.80,  $B(x, \zeta) \cap \text{Ext}(T, \mathbb{R}\langle\zeta\rangle) \cap \text{Zer}(X_{k+1}, \mathbb{R}\langle\zeta\rangle^{k+1})$  is non-empty and contains an element  $z$ . It is clear that  $\lim_\zeta(z) = x$ .

If  $S'$  is bounded over  $\mathbb{R}$  by  $M$  and semi-algebraically connected, then, by Theorem 5.46 (Semi-algebraic triviality), there exists a positive  $t$  in  $\mathbb{R}$  such that  $T_{(0, 2t)}$  is semi-algebraically homeomorphic to  $T_t \times (0, 2t)$ . Thus  $\text{Ext}(T_t, \mathbb{R}\langle\zeta\rangle)$  is semi-algebraically homeomorphic to  $T_\zeta = S$ , which is semi-algebraically connected. Thus  $T_t$  and  $T_{(0, t)}$  are semi-algebraically connected. It follows that

$$S = \overline{T} \cap \text{Zer}(X_{k+1}, \mathbb{R}^{k+1}) = \overline{T} \cap (B(0, M) \times [0, t]) \cap \text{Zer}(X_{k+1}, \mathbb{R}^{k+1})$$

is semi-algebraically connected. □

**Proof of Proposition 12.42:** The proposition follows from

$$\lim_\zeta (\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle\zeta\rangle^k)) = \text{Zer}(Q, \mathbb{R}^k),$$

since  $\text{Zer}(\text{Cr}(Q, \zeta), \mathbb{R}\langle\zeta\rangle^k)$  meets every connected component of

$$\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle\zeta\rangle^k)$$

by Proposition 7.4 and the image of a bounded semi-algebraically connected semi-algebraic set under  $\lim_{\zeta}$  is again semi-algebraically connected by Proposition 12.43.  $\square$

Moreover, the polynomial system  $\text{Cr}(Q, \zeta)$  has good algebraic properties.

**Proposition 12.44.**

- a) *The polynomial system  $\text{Cr}(Q, \zeta)$  is a Gröbner basis for the graded lexicographical ordering with  $X_1 >_{\text{grlex}} \dots >_{\text{grlex}} X_k$ .*
- b) *The set  $\text{Zer}(\text{Cr}(Q, \zeta), \mathbb{C}\langle \zeta \rangle^k)$  is finite.*
- c) *The zeros of the polynomial system  $\text{Cr}(Q, \zeta)$  are simple.*

For the proof of the proposition, we need the following lemma.

**Lemma 12.45.** *The polynomial system*

$$\text{Cr}(Q, 1) = \left\{ G_k(\bar{d}, c), \frac{\partial G_k(\bar{d}, c)}{\partial X_2}, \dots, \frac{\partial G_k(\bar{d}, c)}{\partial X_k} \right\}$$

*has a finite number of zeros in  $\mathbb{C}^k$  all of which are simple.*

**Proof:** Since

$$\frac{\partial G_k(\bar{d}, c)}{\partial X_i} = c^{\bar{d}_1} (\bar{d}_i X_i^{\bar{d}_i - 1} + 2 X_i),$$

for  $i > 1$ , and the zeros of  $\bar{d}_i X_i^{\bar{d}_i - 1} + 2 X_i$  in  $\mathbb{C}$  are simple, the zeros of

$$\left\{ \frac{\partial G_k(\bar{d}, c)}{\partial X_2}, \dots, \frac{\partial G_k(\bar{d}, c)}{\partial X_k} \right\}$$

in  $\mathbb{C}^{k-1}$  are simple and finite in number. A zero of  $\text{Cr}(Q, 1)$  in  $\mathbb{C}^k$  corresponds to a zero  $(x_2, \dots, x_k)$  of

$$\frac{\partial G_k(\bar{d}, c)}{\partial X_2}, \dots, \frac{\partial G_k(\bar{d}, c)}{\partial X_k}$$

in  $\mathbb{C}^{k-1}$  and a zero of  $G_k(\bar{d}, c)(X_1, x_2, \dots, x_k)$  in  $\mathbb{C}$ . Since  $x_i$ ,  $i = 2, \dots, k$ , has norm less than 1 and  $c \leq 1$ ,

$$G_k(\bar{d}, c)(X_1, x_2, \dots, x_k) = c^{\bar{d}_1} X_1^{\bar{d}_1} + a,$$

with  $a$  non-zero, has a finite number of zeros, and all its zeros are simple. This proves the claim.  $\square$

**Proof of Proposition 12.44:** The polynomial system  $\text{Cr}(Q, \zeta)$  is a Gröbner basis for the graded lexicographical ordering according to Proposition 12.3. The set  $\text{Zer}(\text{Cr}(Q, \zeta), \mathbb{C}\langle \zeta \rangle^k)$  is finite according to Corollary 12.9. Consider, for every  $b \neq 0 \in \mathbb{C}$ ,

$$\text{Def}(Q, b) = b G_k(\bar{d}, c) + (1 - b) Q.$$

The polynomial system

$$\text{Cr}(Q, b) = \left\{ \text{Def}(Q, b), \frac{\partial \text{Def}(Q, b)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, b)}{\partial X_k} \right\}$$

is a Gröbner basis for the graded lexicographical ordering according to Proposition 12.3. The set  $\text{Zer}(\text{Cr}(Q, b), \mathbb{C}^k)$  is finite according to Corollary 12.9. We denote by  $A_b$  the finite dimensional vector space

$$A_b = \mathbb{R}[X_1, \dots, X_k] / \text{Ideal}(\text{Cr}(Q, b), \mathbb{R}).$$

Let  $a$  be a separating element of  $\text{Zer}(\text{Cr}(Q, 1), \mathbb{C}^k)$ . According to Lemma 12.45, the zeros of  $\text{Cr}(Q, 1)$  are simple, thus the characteristic polynomial  $\chi_1(T)$  of the linear map  $L_a$  from  $A_1$  to  $A_1$  has only simple roots by Proposition 12.16.

Denoting the characteristic polynomial of the linear map  $L_a$  from  $A_b$  to  $A_b$  by  $\chi_b(T)$ ,

$$B = \{b \in \mathbb{C} \mid b = 0 \text{ or } b \neq 0 \text{ and } \text{Disc}_T(\chi_b(T)) = 0\}$$

is an algebraic subset of  $\mathbb{C}$  which does not contain 1 and is thus finite (see Exercise 1.1). It is clear that  $\zeta \notin \text{Ext}(B, \mathbb{C}\langle\zeta\rangle)$  (see Exercise 1.12). So,  $\text{Disc}_T(\chi_\zeta(T)) \neq 0$ , and by Proposition 4.18, the characteristic polynomial of the linear map  $L_a$  from

$$A_\zeta = \mathbb{R}\langle\zeta\rangle[X_1, \dots, X_k] / \text{Ideal}(\text{Cr}(Q, \zeta), \mathbb{R}\langle\zeta\rangle)$$

to  $A_\zeta$  has only simple zeros. Hence, by Theorem 4.97 (Stickelberger), the zeros of  $\text{Cr}(Q, \zeta)$  are simple. □

**Notation 12.46.** We need to modify slightly the polynomial system

$$\text{Cr}(Q, \zeta) = \left\{ \text{Def}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k} \right\}$$

defined in Notation 12.35 in order to obtain a special Gröbner basis.

Note that defining  $Q_i$ ,  $1 < i \leq k$ , by

$$\frac{\partial \text{Def}(Q, \zeta)}{\partial X_i} = \bar{d}_i \zeta c^{\bar{d}_1} X_i^{\bar{d}_i - 1} + Q_i,$$

we have  $\deg(Q_i) < \bar{d}_i - 1$ ,  $\deg_{X_j}(Q_i) < \bar{d}_j - 1$ ,  $j \neq i$ ,  $1 \leq j \leq k$ , so that  $\text{Cr}(Q, \zeta)$  is nearly a special Gröbner basis. The only properties that are not satisfied are that, defining  $R$  by  $\text{Def}(Q, \zeta) = \zeta c^{\bar{d}_1} X_1^{\bar{d}_1} + R$ , we do not have  $\deg(R) < \bar{d}_1$ , and  $\deg_{X_j}(R) < \bar{d}_j - 1$ ,  $2 \leq j \leq k$ . With  $d = \bar{d}_2 \dots \bar{d}_k$ , we only have to reduce

$$d^2 \zeta^{2k-3} c^{(2k-3)\bar{d}_1} \text{Def}(Q, c)$$

twice modulo each polynomial

$$\frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k}$$

to obtain a polynomial

$$\overline{\text{Def}}(Q, \zeta) = b X_1^{\bar{d}_1} + R_1 \in \mathbb{D}[X_1, \dots, X_k]$$

with  $\deg(R_1) < \bar{d}_i - 1$ ,  $\deg_{X_j}(R_1) < \bar{d}_j - 1$ ,  $j \neq 1$ .

Let

$$\overline{\text{Cr}}(Q, \zeta) = \left\{ \overline{\text{Def}}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k} \right\}.$$

□

It is clear that  $\overline{\text{Cr}}(Q, \zeta)$  is a special Gröbner basis.

Note that  $\text{Cr}(Q, \zeta)$  and  $\overline{\text{Cr}}(Q, \zeta)$  have the same set of zeros.

We are now ready to describe an algorithm giving a point in every connected component of a bounded algebraic set. We simply compute pseudo-critical values and their limits.

*Algorithm 12.16. [Bounded Algebraic Sampling]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $Q(x) \geq 0$  for every  $x \in R^k$  and such that  $\text{Zer}(Q, R^k)$  is contained in  $B(0, 1/c)$ .
- **Output:** a set  $\mathcal{U}$  of real univariate representations of the form

$$(f(T), g(T), \sigma), \text{ with } (f, g) \in D[T]^{k+2}$$

The set of points associated to these univariate representations meets every semi-algebraically connected component of  $\text{Zer}(Q, R^k)$  and contains the set of  $X_1$ -pseudo-critical points on  $\text{Zer}(Q, R^k)$ .

- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degree of  $Q$ .
- **Procedure:**
  - Choose  $(d_1, \dots, d_k)$  such that  $d_1 \geq \dots \geq d_k$ ,  $\deg(Q) \leq d_1$ ,  $t\text{Deg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ . Take as  $\bar{d}_i$  the smallest even number  $> d_i$ ,  $i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ .
  - Compute  $\overline{\text{Cr}}(Q, \zeta)$  (Notation 12.46).
  - Compute the multiplication table  $\mathcal{M}$  of  $\overline{\text{Cr}}(Q, \zeta)$  by Algorithm 12.9 (Special Multiplication Table).
  - Apply the  $\lim_\zeta$  map using Algorithm 12.14 (Limit of Real Bounded Points) with input  $\mathcal{M}$ , and obtain a set  $\mathcal{U}$  of real univariate representations  $v$  with

$$v = (f(T), g(T), \sigma), (f(T), g(T)) \in D[T]^{k+2}.$$

**Proof of correctness:** This follows from Proposition 12.42 and the correctness of Algorithm 12.9 (Special Multiplication Table) and Algorithm 12.14 (Limit of Real Bounded Points). □

**Complexity analysis:** Using the complexity analysis of Algorithm 12.9 (Special Multiplication Table) and Algorithm 12.14 (Limit of Real Bounded Points), the complexity is  $(d_1 \dots d_k)^{O(1)}$  in the ring  $D$ . The polynomials output are of degree  $O(d_1) \dots O(d_k)$  in  $T$ .



When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $Q$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring in the computations of the multiplication table and its output are

$$O(d_1 + \dots + d_{k-1})k d_1(\tau + \nu'),$$

where  $\nu'$  is the bitsize of  $O(d_1 + \dots + d_k)$ , according to the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).

Finally the complexity is  $d^{O(k)}$ , the degree of the univariate representations output are  $O(d)^k$  and the bitsizes of the output are bounded by  $\tau d^{O(k)}$ .  $\square$

*Algorithm 12.17. [Algebraic Sampling]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$ .
- **Output:** a set  $\mathcal{U}$  of real univariate representations of the form

$$(f, g, \sigma), \text{ with } (f, g) \in D[\varepsilon][T]^{k+2}.$$

The set of points associated to these univariate representations meets every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}\langle\varepsilon\rangle^k)$ .

- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degree of  $Q$ .
- **Procedure:**
  - Define

$$R := Q^2 + (\varepsilon(X_1^2 + \dots + X_{k+1}^2) - 1)^2.$$

- Apply Algorithm 12.16 (Bounded Algebraic Sampling) to  $R$ , and obtain a set  $\mathcal{V}$  of real univariate representations  $v$  with

$$v = (f(T), h(T), \sigma), \text{ with } (f, h) \in D[\varepsilon][T]^{k+2}.$$

Define  $\pi(v)$  by  $u$ , with

$$u = (f(T), h_0(T), \dots, h_k(T), \sigma).$$

and  $\mathcal{U} = \pi(\mathcal{V})$ .

**Proof of correctness:** This follows from Proposition 12.34 and the correctness of Algorithm 12.16 (Bounded Algebraic Sampling).  $\square$

**Complexity analysis:** Using the complexity analysis of Algorithm 12.16 (Bounded Algebraic Sampling), and since the degree of  $R$  with respect to  $X_{k+1}$  is 4, the complexity is  $(d_1 \dots d_k)^{O(1)}$  in the ring  $D[\varepsilon]$ . The polynomials output are of degree  $O(d_1) \dots O(d_k)$  in  $T$ . Moreover the degrees with respect to  $\varepsilon$  occurring in the computations of the multiplication table are bounded by

$$O(d_1 + \dots + d_{k-1})k d_k,$$

according to the multiplicity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $Q$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring in the computations of the multiplication table and its output are

$$O(d_1 + \dots + d_{k-1})k d_1(\tau + \nu'),$$

where  $\nu'$  is the bitsize of  $O(d_1 + \dots + d_k)$ , according to the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).

Finally the complexity is  $d^{O(k)}$ , the degree of the univariate representations output in  $T$  and  $\varepsilon$  are  $O(d)^k$  and the bitsizes of the output are bounded by  $d^{O(k)}$ .  $\square$

The following parametrized version of Algorithm 12.16 (Bounded Algebraic Sampling) will be useful in later chapters.

*Algorithm 12.18. [Parametrized Bounded Algebraic Sampling]*

- **Structure:** an ordered integral domain  $D$ .
- **Input:** a polynomial  $Q \in D[Y, X_1, \dots, X_k]$ , such that  $Q(y, x) \geq 0$  for every  $x \in \mathbb{R}^k$ ,  $y \in \mathbb{R}^\ell$ , and for every  $y \in \mathbb{R}^\ell$   $\text{Zer}(Q(y), \mathbb{R}^k)$  is contained in  $B(0, 1/c)$ .
- **Output:** a set  $\mathcal{U}$  of parametrized univariate representations of the form

$$(f, g) \in D[Y, T]^{k+2}.$$

For every  $y \in \mathbb{R}^\ell$ , the set of points associated to these univariate representations meets every semi-algebraically connected component of  $\text{Zer}(Q(y), \mathbb{R}^k)$  and contains the set of  $X_1$ -pseudo-critical points on  $\text{Zer}(Q(y), \mathbb{R}^k)$ .

- **Complexity:**  $(\lambda d^k)^{O(\ell)}$ , where  $d$  is a bound on the degree of  $Q$  with respect to  $X$  and  $\lambda$  is a bound on the degree of  $Q$  with respect to  $Y$ .
- **Procedure:**
  - Choose  $(d_1, \dots, d_k)$  such that  $d_1 \geq \dots \geq d_k$ ,  $\deg(Q) \leq d_1$ ,  $\text{tDeg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ . Take as  $\bar{d}_i$  the smallest even number  $> d$ ,  $i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ .
  - Consider  $\overline{\text{Cr}}(Q, \zeta)$ , using Notation 12.46.
  - Compute the parametrized multiplication table  $\mathcal{M}$  of  $\overline{\text{Cr}}(Q, \zeta)$  by Algorithm 12.10 (Parametrized Special Multiplication Table).
  - Apply Algorithm 12.15 (Parametrized Limit of Bounded Points) with input  $\mathcal{M}$  and  $\zeta$  and obtain a set  $\mathcal{U}$  of parametrized univariate representations  $(v, \sigma)$  with

$$u = f(T), g(T) \in D[Y, T]^{k+2}.$$

**Proof of correctness:** Follows from Proposition 12.42, and the correctness of Algorithm 12.10 (Parametrized Special Multiplication Table) and Algorithm 12.15 (Parametrized Limit of Bounded Points).  $\square$

**Complexity analysis:** Using the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table) and Algorithm 12.15 (Parametrized Limit of Bounded Points), the complexity is  $(d_1, \dots, d_k)^{O(1)}$  in the ring  $D[Y]$ . The polynomials output are of degree  $O(d_1) \dots O(d_k)$  in  $T$  and, if  $\lambda$  is a bound on the total degree in  $Y = (Y_1, \dots, Y_\ell)$  of  $Q$ , of degrees  $\lambda(d_1 \dots d_k)^{O(1)}$  in  $Y$ . Finally, the complexity is  $(\lambda d_1 \dots d_k)^{O(\ell)}$  in the ring  $D$ . The number of elements of  $\mathcal{U}$  is  $O(d_1) \dots O(d_k)$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $Q$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring in the computations of the multiplication table and its output are

$$O(d_1 + \dots + d_{k-1})k d_1(\tau + \nu'),$$

where  $\nu'$  is the bitsize of  $O(\lambda(d_1 + \dots + d_k))^{\ell+1}$ , according to the complexity analysis of 12.10 (Parametrized Special Multiplication Table). Finally the complexity is  $(\lambda d^k)^{O(\ell)}$ , the degree of the univariate representations output are  $O(d)^k$  and the bitsizes of the output are bounded by  $O(k^2 d^2(\tau + \ell \log_2(kd)))$ . □

## 12.7 Triangular Sign Determination

We now give algorithms for sign determination and Thom’s encodings of triangular systems (Definition 11.3). These algorithms have a slightly better complexity than the similar recursive ones in Chapter 11 and will also be easier to generalize to a parametrized setting in later chapters.

We need a notation: let  $u = (f, g) \in K[T]^{k+2}$ ,  $g = (g_0, \dots, g_k)$  be a  $k$ -univariate representation and  $Q \in K[X_1, \dots, X_k]$ . Set

$$Q_u = g_0^e Q\left(\frac{g_k}{g_0}, \dots, \frac{g_k}{g_0}\right), \tag{12.18}$$

where  $e$  is the least even number not less than the degree of  $Q$ .

*Algorithm 12.19.* **[Triangular Sign Determination]**

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $K$ .
- **Input:** a triangular system  $\mathcal{T}$ , and a list  $\mathcal{Q}$  of elements of  $D[X_1, \dots, X_k]$ . Denote by  $Z = \text{Zer}(\mathcal{T}, \mathbb{R}^k)$ .
- **Output:** the set  $\text{SIGN}(\mathcal{Q}, \mathcal{T})$  of sign conditions realized by  $\mathcal{Q}$  on  $Z$ .
- **Complexity:**  $s d'^{O(k)} d^{O(1)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{Q}$ ,  $d'$  is a bound on the degrees on the elements of  $\mathcal{T}$ , and  $d$  is a bound on the degrees of  $Q \in \mathcal{Q}$ .

- **Procedure:** Apply Algorithm 12.16 (Bounded Algebraic Sampling) to  $P = \sum_{A \in \mathcal{T}} A^2$ . Let  $\mathcal{U}$  be the set of real univariate representations output. Keep those  $(u, \sigma) \in \mathcal{U}$ , with  $u = (f, g_0, g_1, \dots, g_k)$ , such that  $P_u$  is zero at the root  $t_\sigma$  of  $f$  with Thom encoding  $\sigma$ , using Algorithm 10.15 (Sign at the Root in a real closed field). For every such real univariate representation  $(u, \sigma)$  and for every  $Q \in \mathcal{Q}$ , compute the sign of  $Q_u$  at  $t_\sigma$ , using Algorithm 10.15 (Sign at the Roots in a real closed field).

**Complexity analysis:** Let  $d'$  be a bound on the degrees of  $P_i$ ,  $d$  a bound on the degrees of  $Q \in \mathcal{Q}$ , and let  $s$  be a bound on the cardinality of  $\mathcal{Q}$ . The number of arithmetic operations in  $\mathbb{D}$  is  $s d'^{O(k)} d^{O(1)}$ , using the complexity analysis of Algorithm 12.16 (Bounded Algebraic Sampling) and Algorithm 10.15 (Sign at the Roots in a real closed field).

When  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{T}$  and  $\mathcal{Q}$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d'^{O(k)} d^{O(1)}$ , using the complexity analysis of Algorithm 12.16 (Bounded Algebraic Sampling) and Algorithm 10.15 (Sign at the Roots in a real closed field).  $\square$

*Algorithm 12.20.* [Triangular Thom Encoding]

- **Structure:** an ordered integral domain  $\mathbb{D}$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a zero-dimensional system of equations  $\mathcal{T}$ .
- **Output:** the list  $\text{Thom}(\mathcal{T})$  of Thom encodings of the roots of  $\mathcal{T}$  (Definition).
- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degrees of  $P \in \mathcal{T}$ .
- **Procedure:** Apply Algorithm 12.19 (Triangular Sign Determination) to  $\mathcal{T}$  and  $\text{Der}(\mathcal{T})$ .

**Complexity analysis:** Since there are  $k d$  polynomials of degrees bounded by  $d$ , in  $\text{Der}(\mathcal{T})$ , the number of arithmetic operations in  $\mathbb{D}$  is  $d^{O(k)}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

When  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{T}$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).  $\square$

*Algorithm 12.21.* [Triangular Comparison of Roots]

- **Structure:** an ordered integral domain  $\mathbb{D}$ , contained in a real closed field  $\mathbb{R}$ .
- **Input:** a Thom encoding  $\mathcal{T}$ ,  $\sigma$  specifying  $z \in \mathbb{R}^{k-1}$ , and two non-zero polynomials  $P$  and  $Q$  in  $\mathbb{D}[X_1, \dots, X_k]$ .
- **Output:** the ordered list of the Thom encodings of the roots of  $P$  and  $Q$  above  $\sigma$  (Definition 11.6).

- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degrees of  $P$ ,  $Q$  and the polynomials in  $\mathcal{T}$ .
- **Procedure:** Apply Algorithm 12.19 (Triangular Sign Determination) to  $\mathcal{T}, P, \text{Der}(\mathcal{T}) \cup \text{Der}(P) \cup \text{Der}(Q)$ , then to  $\mathcal{T}, Q, \text{Der}(\mathcal{T}) \cup \text{Der}(Q) \cup \text{Der}(P)$ . Compare the roots using Proposition 2.28.

**Complexity analysis:** Since there are  $(k + 1)d$  polynomials of degrees bounded by  $d$  in  $\text{Der}(\mathcal{T}) \cup \text{Der}(P) \cup \text{Der}(Q)$ , the number of arithmetic operations in  $D$  is  $d^{O(k)}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{T}$ ,  $P$  and  $Q$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).  $\square$

We can also construct points between two consecutive roots.

*Algorithm 12.22.* **[Triangular Intermediate Points]**

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a Thom encoding  $\mathcal{T}$ ,  $\sigma$  specifying  $z \in \mathbb{R}^{k-1}$ , and two non-zero polynomials  $P$  and  $Q$  in  $D[X_1, \dots, X_k]$ .
- **Output:** Thom encodings specifying values  $y$  intersecting intervals between two consecutive roots of  $P(z, X_k)$  and  $Q(z, X_k)$ .
- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degrees of  $P$ ,  $Q$  and the polynomials in  $\mathcal{T}$ .
- **Procedure:** Compute the Thom encodings of the roots of

$$\partial(PQ)/\partial X_k(z, X_k)$$

above  $\mathcal{T}$ ,  $\sigma$  using Algorithm 12.20 (Triangular Thom Encoding) and compare them to the roots of  $P$  and  $Q$  above  $\sigma$  using Algorithm 12.21 (Triangular Comparison of Roots). Keep one intermediate point between two consecutive roots of  $PQ$ .

**Complexity analysis:** Then the degree of  $\partial(PQ)/\partial X_k(z, X_k)$  is  $O(d)$ . Using the complexity of Algorithms 12.20 (Triangular Thom Encoding) and Algorithm 12.21 (Triangular Comparison of Roots), the complexity is  $d^{O(k)}$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{T}$ ,  $P$  and  $Q$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ , using the complexity analysis of Algorithms 12.20 (Triangular Thom Encoding) and Algorithm 12.21 (Triangular Comparison of Roots).  $\square$

Finally we can compute sample points on a line.

*Algorithm 12.23. [Triangular Sample Points]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a Thom encoding  $\mathcal{T}$ ,  $\sigma$  specifying  $z \in R^{k-1}$ , and a family of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** an ordered list  $L$  of Thom encodings specifying the roots in  $R$  of the non-zero polynomials  $P(z, X_k)$ ,  $P \in \mathcal{P}$ , an element between two such consecutive roots, an element of  $R$  smaller than all these roots, and an element of  $R$  greater than all these roots. Moreover  $(\tau_1)$  appears before  $(\tau_2)$  in  $L$  if and only if  $x_k(\tau_1) \leq x_k(\tau_2)$ . The sign of  $Q(z, x_k(\tau))$  is also output for every  $Q \in \mathcal{P}$ ,  $\tau \in L$ .
- **Complexity:**  $s^2 d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{T}$  and  $\mathcal{Q}$ .
- **Procedure:** Characterize the roots of the polynomials in  $R$  using Algorithm 12.20 (Triangular Thom Encoding). Compare these roots using Algorithm 12.21 (Triangular Comparison of Roots) for every pair of polynomials in  $\mathcal{P}$ . Characterize a point in each interval between the roots by Algorithm 12.22 (Triangular Intermediate Points). Use Proposition 10.5 to find an element of  $R$  smaller and bigger than any root of any polynomial in  $\mathcal{P}$  above  $z$ .

**Complexity analysis:** Using the complexity analyses of Algorithm 12.20 (Triangular Thom Encoding), Algorithm 12.21 (Triangular Comparison of Roots) and Algorithm 12.22 (Triangular Intermediate Points), the complexity is  $s^2 d^{O(k)}$ .

When  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $\mathcal{T}$  and  $\mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ , using the complexity analysis of Algorithm 12.20 (Triangular Thom Encoding), Algorithm 12.21 (Triangular Comparison of Roots) and Algorithm 12.22 (Triangular Intermediate Points).  $\square$

## 12.8 Computing the Euler-Poincaré Characteristic of an Algebraic Set

In this section we first describe an algorithm for computing the Euler-Poincaré characteristic of an algebraic set. The complexity of this algorithm is asymptotically the same as that of Algorithm 12.16 (Bounded Algebraic Sampling) for computing sample points in every connected component of a bounded algebraic set described in the Section 12.6.

We first describe an algorithm for computing the Euler-Poincaré characteristic of a bounded algebraic set and then use this algorithm for computing the Euler-Poincaré characteristic of a general algebraic set.

From now on we consider a polynomial  $Q \in D[X_1, \dots, X_k]$ , where  $D$  is a ring contained in the real closed field  $R$ , satisfying  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$  for some  $0 < c \leq 1, c \in D$ . Let  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$  with  $\bar{d}_i$  even and  $\bar{d}_i > d_i$ , where  $d_i$  is the total degree of  $Q^2$  in  $X_i$ .

**Notation 12.47.** We denote

$$\begin{aligned} G_k(\bar{d}, c) &= c^{\bar{d}_1}(X_1^{\bar{d}_1} + \dots + X_k^{\bar{d}_k} + X_2^2 + \dots + X_k^2) - (2k - 1), \\ \text{Def}(Q^2, \zeta) &= \zeta G_k(\bar{d}, c) + (1 - \zeta)Q^2, \\ \text{Def}_+(Q^2, \zeta) &= \text{Def}(Q^2, \zeta) + X_{k+1}^2, \\ \text{Cr}(Q^2, \zeta) &= \left\{ \text{Def}(Q^2, \zeta), \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_k} \right\}, \\ \text{Cr}_+(Q^2, \zeta) &= \left\{ \text{Def}_+(Q^2, \zeta), \frac{\partial \text{Def}_+(Q^2, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}_+(Q^2, \zeta)}{\partial X_k}, 2X_{k+1} \right\}, \\ \overline{\text{Cr}}(Q^2, \zeta) &= \left\{ \overline{\text{Def}}(Q^2, \zeta), \frac{\partial \overline{\text{Def}}(Q^2, \zeta)}{\partial X_2}, \dots, \frac{\partial \overline{\text{Def}}(Q^2, \zeta)}{\partial X_k} \right\}, \\ \overline{\text{Cr}}_+(Q^2, \zeta) &= \left\{ \overline{\text{Def}}_+(Q^2, \zeta), \frac{\partial \overline{\text{Def}}_+(Q^2, \zeta)}{\partial X_2}, \dots, \frac{\partial \overline{\text{Def}}_+(Q^2, \zeta)}{\partial X_k} \right\} \end{aligned}$$

where,  $\overline{\text{Def}}_+(Q^2, \zeta)$  is obtained from  $\text{Def}_+(Q^2, \zeta)$  as in Notation 12.46.  $\square$

It is clear that  $\overline{\text{Cr}}(Q^2, \zeta)$  as well as  $\overline{\text{Cr}}_+(Q^2, \zeta)$  are both special Gröbner bases. Note that  $\text{Cr}(Q^2, \zeta)$  and  $\overline{\text{Cr}}(Q^2, \zeta)$  (resp.  $\text{Cr}_+(Q^2, \zeta)$  and  $\overline{\text{Cr}}_+(Q^2, \zeta)$ ) have the same set of zeros.

*Algorithm 12.24.* **[Euler-Poincaré Characteristic of a Bounded Algebraic Set]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$  for which  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$ .
- **Output:** the Euler-Poincaré characteristic  $\chi(\text{Zer}(Q, R^k))$ .
- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degree of  $Q$ .
- **Procedure:**
  - Choose  $(d_1, \dots, d_k)$  such that  $d_1 \geq \dots \geq d_k$ ,  $\deg(Q^2) \leq d_1$ , and  $\text{tDeg}_{X_i}(Q^2) \leq d_i$ , for  $i = 2, \dots, k$ . Take  $\bar{d}_i$  the smallest even number  $> d_i, i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ .
  - Consider  $\overline{\text{Cr}}(Q^2, \zeta)$  and  $\overline{\text{Cr}}_+(Q^2, \zeta)$ , using Notation 12.47.
  - Compute the multiplication tables  $\mathcal{M}$  and  $\mathcal{M}_+$  of  $\overline{\text{Cr}}(Q^2, \zeta)$  and  $\overline{\text{Cr}}_+(Q^2, \zeta)$  using Algorithm 12.10 (Parametrized Special Multiplication Table), with parameter  $\zeta$ .
  - Compute the characteristic polynomial of the matrices

$$H_1 = \left[ \frac{\partial^2 \text{Def}(Q^2, \zeta)}{\partial X_i \partial X_j} \right]_{2 \leq i, j \leq k}$$

and,

$$H_2 = \left[ \frac{\partial^2 \text{Def}_+(Q^2, \zeta)}{\partial X_i \partial X_j} \right]_{2 \leq i, j \leq k+1}$$

using Algorithm 8.17 (Characteristic Polynomial).

- Compute the signature  $\text{Sign}(H_1)$  (resp.  $\text{Sign}(H_2)$ ), of the matrix  $H_1$  (resp.  $H_2$ ) at the real roots of  $\text{Cr}(Q^2)$  (resp.  $\text{Cr}_+(Q^2)$ ) using Algorithm 12.8 (Multivariate Sign Determination) with input the list of coefficients of the characteristic polynomial of  $H_1$  (resp.  $H_2$ ) and the multiplication table  $\mathcal{M}$  (resp.  $\mathcal{M}_+$ ) for the zero-dimensional system  $\overline{\text{Cr}}(Q^2)$  (resp.  $\overline{\text{Cr}}_+(Q^2)$ ), to determine the signs of the coefficients of characteristic polynomials of  $H_1$  (resp.  $H_2$ ) at the real roots of the corresponding system.
- For  $i$  from 0 to  $k-1$  let,

$$\ell_i := \#\{x \in \text{Zer}(\text{Cr}(Q^2), C\langle \zeta \rangle^k) \mid k-1 + \text{Sign}(H_1(x))/2 = i\}.$$

- For  $i$  from 0 to  $k$ , let

$$m_i := \#\{x \in \text{Zer}(\text{Cr}_+(Q^2), C\langle \zeta \rangle^{k+1}) \mid k + \text{Sign}(H_2(x))/2 = i\}.$$

- Output

$$\chi(\text{Zer}(Q, \mathbb{R}^k)) = \frac{1}{2} \left( \sum_{i=0}^{k-1} (-1)^{k-1-i} \ell_i + \sum_{i=0}^k (-1)^{k-i} m_i \right).$$

We need the following lemma.

**Lemma 12.48.** *The Hessian matrices  $H_1$  (resp.  $H_2$ ) are non-singular at the points of  $\text{Zer}(\text{Cr}(Q^2, \zeta), C\langle \zeta \rangle^k)$  (resp.  $\text{Zer}(\text{Cr}_+(Q^2, \zeta), C\langle \zeta \rangle^{k+1})$ ).*

**Proof:** Let

$$\begin{aligned} \text{Def}(Q^2, \lambda, \mu) &= \lambda Q^2 + \mu G(\bar{d}, c), \\ \text{Def}^h(Q^2, \lambda, \mu) &= \lambda (Q^2)^h + \mu G(\bar{d}, c)^h \end{aligned}$$

being its homogenization in degree  $\bar{d}_1$ .

Moreover, let

$$H_1(\lambda, \mu) = \left[ \frac{\partial^2 \text{Def}(Q^2, \lambda, \mu)}{\partial X_i \partial X_j} \right]_{2 \leq i, j \leq k}$$

be the corresponding Hessian matrix and  $\text{Cr}^h(Q^2, \lambda, \mu, )$  the corresponding system of equations for the polynomial  $\text{Def}^h(Q^2, \lambda, \mu, )$ .

Now,  $H_1(0, 1)$  is the diagonal matrix with entries

$$\bar{d}_i(\bar{d}_i - 1) X_i^{\bar{d}_i - 2} + 2, \quad 2 \leq i, j \leq k.$$



Also, if  $(x_1, \dots, x_k) \in \text{Zer}(\text{Cr}(Q^2, 0, 1), \mathbb{C}^k)$ , then each  $x_i$ , for  $i = 2, \dots, k$  has norm less than 1 (see proof of Lemma 12.45). Also,  $\text{Zer}(\text{Cr}^h(Q^2, 0, 1), \mathbb{P}_k(\mathbb{C}))$  has no points at infinity. Hence, it is clear that  $\det^h(H_1(0, 1)) \neq 0$  at any point of  $\text{Zer}(\text{Cr}^h(Q^2, 0, 1), \mathbb{P}_k(\mathbb{C}))$ .

Thus, the set  $D$  of  $(\lambda : \mu) \in \mathbb{P}_1(\mathbb{C})$  such that  $\det^h(H_1(\lambda, \mu)) = 0$  at a point of  $\text{Zer}(\text{Cr}^h(Q^2, \lambda, \mu), \mathbb{P}_k(\mathbb{C}))$  does not contain  $(0 : 1)$ .

Moreover,  $D$  is the projection on  $\mathbb{P}_1(\mathbb{C})$  of an algebraic subset of  $\mathbb{P}_k(\mathbb{C}) \times \mathbb{P}_1(\mathbb{C})$  and is thus algebraic by Theorem 4.102. Since,  $D$  does not contain the point  $(0 : 1)$  it is a finite subset of  $\mathbb{P}_1(\mathbb{C})$  by Lemma 4.101. Hence, the set of  $t \in \mathbb{C}$  such that  $\det^h(H_1(1 - t, t)) = 0$  at a point of  $\text{Zer}(\text{Cr}(Q^2, 1 - t, t), \mathbb{C})$  is finite and its extension to  $\mathbb{C}(\zeta)$  is a finite number of elements of  $\mathbb{C}$  which does not contain  $\zeta$ .

This claim now follows for  $H_1$  and the proof is identical for  $H_2$ . □

**Proof of correctness of Algorithm 12.24:** It follows from Proposition 12.38 and Lemma 12.40 that the algebraic sets  $\text{Zer}(\text{Def}(Q^2, -\zeta), \mathbb{R}(\zeta)^k)$  and  $\text{Zer}(\text{Def}_+(Q^2, \zeta), \mathbb{R}(\zeta)^{k+1})$  are non-singular algebraic hypersurfaces bounded over  $\mathbb{R}$ .

Moreover, by Proposition 12.44, the zeros of the polynomial system  $\text{Cr}(Q^2, \zeta)$  are simple. The same holds for  $\text{Cr}_+(Q^2, \zeta)$ .

It follows from Lemma 12.48 that the projection map onto the  $X_1$  coordinate has non-degenerate critical points on the hypersurfaces  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}(\zeta)^k)$  and  $\text{Zer}(\text{Def}_+(Q^2, \zeta), \mathbb{R}(\zeta)^{k+1})$ .

Now, the algebraic set  $\text{Zer}(\text{Def}_+(Q^2, \zeta), \mathbb{R}(\zeta)^{k+1})$  is semi-algebraically homeomorphic to two copies of the set,  $S$  defined by  $\text{Def}(Q^2, \zeta) \leq 0$ , glued along  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}(\zeta)^k)$ .

It follows from Equation 6.36 that,

$$\chi(\text{Zer}(\text{Def}_+(Q^2, \zeta), \mathbb{R}(\zeta)^{k+1})) = 2\chi(S) - \chi(\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}(\zeta)^k)).$$

We now claim that the closed and bounded set  $S$  has the same homotopy type (and hence the same Euler-Poincaré characteristic) as  $\text{Zer}(Q, \mathbb{R}(\zeta)^k)$ .

We replace  $\zeta$  in the definition of the set  $S$  by a new variable  $T$ , and consider the set  $K \subset \mathbb{R}^{k+1}$  defined by  $\{(x, t) \in \mathbb{R}^{k+1} \mid \text{Def}(Q^2, T) \leq 0\}$  and let for  $b > 0$ ,  $K_b = \{x \mid (x, b) \in K\}$ . Note that  $\text{Ext}(K, \mathbb{R}(\zeta))_\zeta = S$ . Let  $\pi_X$  (resp.  $\pi_T$ ) denote the projection map onto the  $X$  (resp.  $T$ ) coordinates.

Clearly,  $\text{Zer}(Q, \mathbb{R}^k) \subset K_b$  for all  $b > 0$ . By Theorem 5.46 (Semi-algebraic triviality), for all small enough  $b > 0$ , there exists a semi-algebraic homeomorphism,  $\phi: K_b \times (0, b] \rightarrow K \cap \pi_T^{-1}((0, b])$ , such that  $\pi_T(\phi(x, s)) = s$  and  $\phi(\text{Zer}(Q, \mathbb{R}^k), s) = \text{Zer}(Q, \mathbb{R}^k)$  for all  $s \in (0, b]$ .

Let  $G: K_b \times [0, b] \rightarrow K_b$  be the map defined by  $G(x, s) = \pi_X(\phi(x, s))$  for  $s > 0$  and  $G(x, 0) = \lim_{s \rightarrow 0^+} \pi_X(\phi(x, s))$ . Let  $g: K_b \rightarrow K_b$  be the map  $G(x, 0)$  and  $i: \text{Zer}(Q, \mathbb{R}^k) \rightarrow K_b$  the inclusion map. Using the homotopy  $G$ , we see that  $i \circ g \sim \text{Id}_{K_b}$ , and  $g \circ i \sim \text{Id}_{\text{Zer}(Q, \mathbb{R}^k)}$ , which shows that  $\text{Zer}(Q, \mathbb{R}^k)$  is homotopy equivalent to  $K_b$  for all small enough  $b > 0$ . Now, specialize  $b$  to  $\zeta$ .

Finally, the correctness of the computations of

$$\chi(\text{Zer}(\text{Def}_+(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})), \chi(\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k))$$

is a consequence of Lemma 7.25, using Tarski-Seidenberg principle (Theorem 2.80).  $\square$

**Complexity analysis of Algorithm 12.24:** The complexity of the algorithm is  $d^{O(k)}$ , according to the complexity analysis of Algorithm 12.10 (Special Multiplication Table), Algorithm 8.17 (Characteristic polynomial), Algorithm 12.8 (Multivariate Sign Determination).

When  $D = \mathbb{Z}$  and the bitsizes of the coefficients of  $Q$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ .  $\square$

*Algorithm 12.25. [Euler-Poincaré Characteristic of an Algebraic Set]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$ .
- **Output:** the Euler-Poincaré characteristic  $\chi(\text{Zer}(Q, \mathbb{R}^k))$ .
- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degree of  $Q$ .
- **Procedure:**
  - Define

$$\begin{aligned} Q_1 &= Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2) - 1)^2, \\ Q_2 &= Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_{k+1}^2) - 1)^2, \end{aligned}$$

- Using Algorithm 12.24 (Euler-Poincaré Characteristic of a Bounded Algebraic Set) compute  $\chi(\text{Zer}(Q_1, \mathbb{R}\langle \varepsilon \rangle^k))$  and  $\chi(\text{Zer}(Q_2, \mathbb{R}\langle \varepsilon \rangle^{k+1}))$ .
- Output,

$$\chi(\text{Zer}(Q, \mathbb{R}^k)) = (\chi(\text{Zer}(Q_2, \mathbb{R}\langle \varepsilon \rangle^{k+1})) - \chi(\text{Zer}(Q_1, \mathbb{R}\langle \varepsilon \rangle^k)))/2.$$

**Proof of correctness:** It is clear that the algebraic sets  $\text{Zer}(Q_1, \mathbb{R}\langle \varepsilon \rangle^k)$  and  $\text{Zer}(Q_2, \mathbb{R}\langle \varepsilon \rangle^{k+1})$  are bounded over  $\mathbb{R}\langle \varepsilon \rangle$  and hence we can apply Algorithm 12.24 to compute their Euler-Poincaré characteristics.

Moreover,  $\text{Zer}(Q_2, \mathbb{R}\langle \varepsilon \rangle^{k+1})$  is semi-algebraically homeomorphic to two copies of  $\text{Zer}(Q, \mathbb{R}\langle \varepsilon \rangle^k) \cap B(0, 1/\varepsilon)$  glued along the algebraic set  $\text{Zer}(Q_1, \mathbb{R}\langle \varepsilon \rangle^k)$ . Hence, using Equation 6.36 we obtain that,

$$\chi(\text{Zer}(Q, \mathbb{R}\langle \varepsilon \rangle^k) \cap B(0, 1/\varepsilon)) = (\chi(\text{Zer}(Q_2, \mathbb{R}\langle \varepsilon \rangle^{k+1})) + \chi(\text{Zer}(Q_1, \mathbb{R}\langle \varepsilon \rangle^k)))/2.$$

The correctness of the algorithm now follows from the fact that,

$$\chi(\text{Zer}(Q, \mathbb{R}^k)) = \chi(\text{Zer}(Q, \mathbb{R}\langle \varepsilon \rangle^k) \cap B(0, 1/\varepsilon)),$$

since  $\text{Zer}(Q, \mathbb{R}\langle \varepsilon \rangle^k)$  and  $\text{Zer}(Q, \mathbb{R}\langle \varepsilon \rangle^k) \cap B(0, 1/\varepsilon)$  are semi-algebraically homeomorphic.  $\square$

**Complexity Analysis:** The complexity of the algorithm is  $d^{O(k)}$  using the complexity analysis of Algorithm 12.24, and the bound  $d^{O(k)}$  on the degree in  $\varepsilon, \zeta$  obtained in the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).

When  $D = \mathbb{Z}$  and the bitsizes of the coefficients of  $Q$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $d^{O(k)}$ .  $\square$

## 12.9 Bibliographical Notes

Gröbner basis have been introduced and studied by B. Buchberger [32, 33]. They are a major tool in computational algebraic geometry and polynomial system solving.

The idea of representing solutions of polynomial systems by Rational Univariate Representation seems to appear first in the work of Kronecker [100]. The use of these representations in computer algebra starts with [4, 134] and is now commonly used. The first single exponential algorithms in time  $d^{O(k)}$  for finding a point in every connected component of an algebraic set can be found in [37, 133]. The notion of pseudo-critical point is introduced in [75]. The notion of well-separating element and the limit process comes from [137]. The algorithm for computing the Euler-Poincaré characteristic for algebraic sets appears in [11].

---

## Existential Theory of the Reals

The **decision problem for the existential theory of the reals** is to decide the truth or falsity of a sentence  $(\exists X_1) \dots (\exists X_k) F(X_1, \dots, X_k)$ , where  $F(X_1, \dots, X_k)$  is a quantifier free formula in the language of ordered fields with coefficients in a real closed field  $\mathbb{R}$ . This problem is equivalent to deciding whether or not a given semi-algebraic set is empty. It is a special case of the general decision problem seen in Chapter 11.

When done by the Cylindrical Decomposition Algorithm of Chapter 11, deciding existential properties of the reals has complexity doubly exponential in  $k$ , the number of variables. But the existential theory of the reals has a special logical structure, since the sentence to decide has a single block of existential quantifiers. We take advantage of this special structure to find an algorithm which is singly exponential in  $k$ .

Our method for solving the existential theory of the reals is to compute the set of realizable sign conditions of the set of polynomials  $\mathcal{P}$  appearing in the quantifier free formula  $F$ . We have already seen in Proposition 7.35 that the set of realizable sign condition of  $\mathcal{P}$  is polynomial in the degree  $d$  and the number  $s$  of polynomials and singly exponential in the number of variables  $k$ . The proof of Proposition 7.35 used Mayer-Vietoris sequence and Theorem 7.23. Our technique here will be quite different, though the main ideas are inspired by the critical point method already used in Chapter 7 and Chapter 12.

In Section 13.1, we describe an algorithm for computing the set of realizable sign conditions, as well as sample points in their realizations, whose complexity is polynomial in  $s$  and  $d$  and singly exponential in  $k$ . This algorithm uses pseudo-critical points introduced in Chapter 12 and additional techniques for achieving general position by infinitesimal perturbations.

In Section 13.1, we describe some applications of the preceding results related to bounding the size of a ball meeting every semi-algebraically connected component of the realization of every realizable sign condition, as well as to certain real and complex decision problems.

In Section 13.3, we describe an algorithm for computing sample points in realizations of realizable sign conditions on an algebraic set taking advantage of the (possibly low) dimension of the algebraic set.

Finally, in Section 13.4 we describe a method for computing the Euler-Poincaré characteristic of all possible sign conditions defined by a family of polynomials.

### 13.1 Finding Realizable Sign Conditions

In this section, let  $\mathcal{P} = \{\mathcal{P}_1, \dots, \mathcal{P}_s\} \subset \mathbb{R}[X_1, \dots, X_k]$ . Recall that we denote by  $\text{SIGN}(\mathcal{P}) \subset \{0, 1, -1\}^{\mathcal{P}}$  the set of all realizable sign conditions for  $\mathcal{P}$  (see Notation 7.29). We are now going to present an algorithm which computes  $\text{SIGN}(\mathcal{P})$ .

We first prove that we can reduce the problem of computing a set of sample points meeting the realizations of every realizable sign conditions of a family of polynomials to the problem already considered in Chapter 12, namely finding points in every semi-algebraically connected component of certain algebraic sets.

**Proposition 13.1.** *Let  $D \subset \mathbb{R}^k$  be a non-empty semi-algebraically connected component of a basic closed semi-algebraic set defined by*

$$P_1 = \dots = P_\ell = 0, P_{\ell+1} \geq 0, \dots, P_s \geq 0.$$

*There exists an algebraic set  $W$  defined by equations*

$$P_1 = \dots = P_\ell = P_{i_1} = \dots = P_{i_m} = 0,$$

*(with  $\{i_1, \dots, i_m\} \subset \{\ell + 1, \dots, s\}$ ) such that a semi-algebraically connected component  $D'$  of  $W$  is contained in  $D$ .*

**Proof:** Consider a maximal set of polynomials

$$\{P_1, \dots, P_\ell, P_{i_1}, \dots, P_{i_m}\},$$

where

$$m = 0 \text{ or } \ell < i_1 < \dots < i_m \leq s,$$

with the property that there exists a point  $p \in D$  where

$$P_1 = \dots = P_\ell = P_{i_1} = \dots = P_{i_m} = 0.$$

Consider the semi-algebraically connected component  $D'$  of the algebraic set defined by

$$P_1 = \dots = P_\ell = P_{i_1} = \dots = P_{i_m} = 0,$$

which contains  $p$ . We claim that  $D' \subset D$ . Suppose that there exists a point  $q \in D'$  such that  $q \notin D$ . Then by Proposition 5.23, there exists a semi-algebraic path  $\gamma: [0, 1] \rightarrow D'$  joining  $p$  to  $q$  in  $D'$ . Denote by  $q'$  the first point of the path  $\gamma$  on the boundary of  $D$ . More precisely, note that

$$A = \{t \in [0, 1] \mid \gamma([0, t]) \subset D\}$$

is a closed semi-algebraic subset of  $[0, 1]$  which does not contain 1. Thus  $A$  is the union of a finite number of closed intervals

$$A = [0, b_1] \cup \dots \cup [a_\ell, b_\ell].$$

Take  $q' = \gamma(b_1)$ . At least one of the polynomials, say  $P_j$ ,  $j \notin \{1, \dots, \ell, i_1, \dots, i_m\}$  must be 0 at  $q'$ . This violates the maximality of the set

$$\{P_1, \dots, P_\ell, P_{i_1}, \dots, P_{i_m}\}.$$

It is clear that if  $D$  is bounded,  $D'$  is bounded. □

**Proposition 13.2.** *Let  $D \subset \mathbb{R}^k$  be a non-empty semi-algebraically connected component of a semi-algebraic set defined by*

$$P_1 = \dots = P_\ell = 0, P_{\ell+1} > 0, \dots, P_s > 0.$$

*There exists an algebraic set  $W \subset \mathbb{R}\langle\varepsilon\rangle^k$  defined by equations*

$$P_1 = \dots = P_\ell = 0, P_{i_1} = \dots = P_{i_m} = \varepsilon$$

*(with  $\{i_1, \dots, i_m\} \subset \{\ell + 1, \dots, s\}$ ) such that there exists a semi-algebraically connected component  $D'$  of  $W$  which is contained in  $\text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$ .*

**Proof:** Consider two points  $x$  and  $y$  in  $D$ . By Proposition 5.23, there is a semi-algebraic path  $\gamma$  from  $x$  to  $y$  inside  $D$ . Since  $\gamma$  is closed and bounded, the semi-algebraic and continuous function  $\min_{\ell+1 \leq i \leq s} (P_i)$  has a strictly positive minimum on  $\gamma$ . The extension of the path  $\gamma$  to  $\mathbb{R}\langle\varepsilon\rangle$  is thus entirely contained inside the subset  $S$  of  $\mathbb{R}\langle\varepsilon\rangle^k$  defined by

$$P_1 = \dots = P_\ell = 0, P_{\ell+1} - \varepsilon \geq 0, \dots, P_s - \varepsilon \geq 0.$$

Thus, there is one non-empty semi-algebraically connected component  $\bar{D}$  of  $S$  containing  $D$ . Applying Proposition 13.1 to  $\bar{D}$  and  $S$ , we get a semi-algebraically connected component  $D'$  of some

$$P_1 = \dots = P_\ell = 0, P_{i_1} = \dots = P_{i_m} = \varepsilon,$$

contained in  $\bar{D}$ . Then  $D' \subset \text{Ext}(D, \mathbb{R}\langle\varepsilon\rangle)$ . □

*Remark 13.3.* Proposition 13.2, Algorithm 12.16 (Bounded Algebraic Sampling), and Algorithm 11.20 (Removal of Infinitesimals) provide an algorithm outputting a set of points meeting every semi-algebraically connected component of the realization of a realizable sign condition of a family  $\mathcal{P}$  of  $s$  polynomials on a bounded algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  with complexity  $2^s d^{O(k)}$  (where  $d$  is a bound on the degree of  $Q$  and the  $P \in \mathcal{P}$ ), considering all possible subsets of  $\mathcal{P}$ . Note that this algorithm does not involve polynomials of degree doubly exponential in  $k$ , in contrast to Algorithm 11.2 (Cylindrical Decomposition). □

**Exercise 13.1.**

- a) Describe precisely the algorithm outlined in the preceding remark and prove its complexity.

b) Describe an algorithm with the same complexity without the hypothesis that  $\text{Zer}(Q, \mathbb{R}^k)$  is bounded.

When  $s$  is bigger than the dimension  $k$  of the ambient space, the algorithm proposed in the preceding remark does not give a satisfactory complexity bound, since the complexity is exponential in  $s$ . Reduction to general position, using infinitesimal deformations, will be the key for a better complexity result.

Let us define precisely the notion of general position that we consider. Let  $\mathcal{P}^* = \{\mathcal{P}_1^*, \dots, \mathcal{P}_s^*\}$ , where for every  $i = 1, \dots, s$ ,  $\mathcal{P}_i^* \subset \mathbb{R}[X_1, \dots, X_k]$  is finite, and such that two distinct elements of  $\mathcal{P}_i^*$  have no common zeros in  $\mathbb{R}^k$ . The family  $\mathcal{P}^*$  is in  **$\ell$ -general position** with respect to  $Q \in \mathbb{R}[X_1, \dots, X_k]$  in  $\mathbb{R}^k$  if no  $\ell + 1$  polynomials belonging to different  $\mathcal{P}_i^*$  have a zero in common with  $Q$  in  $\mathbb{R}^k$ .

The family  $\mathcal{P}^*$  is in **strong  $\ell$ -general position** with respect to  $Q \in \mathbb{R}[X_1, \dots, X_k]$  in  $\mathbb{R}^k$  if moreover any  $\ell$  polynomials belonging to different  $\mathcal{P}_i^*$  have at most a finite number of zeros in common with  $Q$  in  $\mathbb{R}^k$ .

When  $Q = 0$ , we simply say that  $\mathcal{P}^* \subset \mathbb{R}[X_1, \dots, X_k]$  is in  $\ell$ -general position (resp. strong  $\ell$ -general position) in  $\mathbb{R}^k$ .

We also need the notion of a family of homogeneous polynomials in general position in  $\mathbb{P}_k(\mathbb{C})$ . The reason for considering common zeros in  $\mathbb{P}_k(\mathbb{C})$  is that we are going to use in our proofs the fact that, in the context of complex projective geometry, the projection of an algebraic set is algebraic. This was proved in Theorem 4.102.

Let  $\mathcal{P}^* = \{\mathcal{P}_1^*, \dots, \mathcal{P}_s^*\}$  where for every  $i = 1, \dots, s$ ,  $\mathcal{P}_i^* \in \mathbb{R}[X_0, X_1, \dots, X_k]$  is homogeneous. The family  $\mathcal{P}^*$  is in  **$\ell$ -general position** with respect to a homogeneous polynomial  $Q^h \in \mathbb{R}[X_0, X_1, \dots, X_k]$  in  $\mathbb{P}_k(\mathbb{C})$  if no more than  $\ell$  polynomials in  $\mathcal{P}_i^*$  have a zero in common with  $Q^h$  in  $\mathbb{P}_k(\mathbb{C})$ .

We first give an example of a finite family of polynomials in general position and then explain how to perturb a finite set of polynomials to get a family in strong general position.

**Notation 13.4.** Define

$$H_k(d, i) = 1 + \sum_{1 \leq j \leq k} i^j X_j^d,$$

$$H_k^h(d, i) = X_0^d + \sum_{1 \leq j \leq k} i^j X_j^d.$$

Note that when  $d$  is even,  $H_k(d, i)(x) > 0$  for every  $x \in \mathbb{R}^k$ . □

**Lemma 13.5.** *For any positive integer  $d$ , the polynomials  $H_k^h(d, i)$ ,  $0 \leq i \leq s$ , are in  $k$ -general position in  $\mathbb{P}_k(\mathbb{C})$ .*

**Proof:** Take  $P(T, X_0, \dots, X_k) = X_0^d + \sum_{1 \leq j \leq k} T^j X_j^d$ . If  $k + 1$  of the  $H_k^h(d, i)$  had a common zero  $\bar{x}$  in  $\mathbb{P}_k(\mathbb{C})$ , substituting homogeneous coordinates of this common zero in  $P$  would give a non-zero univariate polynomial in  $T$  of degree at most  $k$  with  $k + 1$  distinct roots, which is impossible. □

Consider three variables  $\varepsilon, \delta, \gamma$  and  $R\langle\varepsilon, \delta, \gamma\rangle$ . Note that  $\varepsilon, \delta, \gamma$  are three infinitesimal quantities in  $R\langle\varepsilon, \delta, \gamma\rangle$  with  $\varepsilon > \delta > \gamma > 0$ . The reason for using these three infinitesimal quantities is the following. The variable  $\varepsilon$  is used to get bounded sets, the variables  $\delta, \gamma$  are used to reach general position, and describe sets which are closely related to realizations of sign conditions on the original family.

Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset R[X_1, \dots, X_k]$  be polynomials of degree bounded by  $d$ . With  $d' > d$ , let  $\mathcal{P}^*$  be the family  $\{P_1^*, \dots, P_s^*\}$  with

$$P_i^* = \{(1 - \delta) P_i + \delta H_k(d', i), (1 - \delta) P_i - \delta H_k(d', i), \\ (1 - \delta) P_i + \delta \gamma H_k(d', i), (1 - \delta) P_i - \delta \gamma H_k(d', i)\}.$$

We prove

**Proposition 13.6.** *The family  $\mathcal{P}^*$  is in strong  $k$ -general position in  $R\langle\varepsilon, \delta, \gamma\rangle^k$ .*

**Proof:** For  $P_i \in \mathcal{P}$  we write

$$P_i^h = X_0^{d'} P_i \left( \frac{X_1}{X_0}, \dots, \frac{X_k}{X_0} \right).$$

Consider

$$\bar{P}_i^*(\lambda, \mu) = \{\lambda P_i^h + \mu H_k^h(d', i), \lambda P_i^h - \mu H_k^h(d', i), \\ \lambda P_i^h + \mu \gamma H_k^h(d', i), \lambda P_i^h - \mu \gamma H_k^h(d', i)\}.$$

Let  $I = \{i_1, \dots, i_{k+1}\}$ , and  $Q_{i_j} \in \bar{P}_{i_j}^*$ ,  $j = 1, \dots, k + 1$ . The set  $D_I$  of  $(\lambda: \mu) \in \mathbb{P}_1(C\langle\gamma\rangle)$  such that

$$Q_{i_1}^h(\lambda, \mu), \dots, Q_{i_{k+1}}^h(\lambda, \mu)$$

have a common zero is the projection on  $\mathbb{P}_1(C\langle\gamma\rangle)$  of an algebraic subset of  $\mathbb{P}_k(C\langle\gamma\rangle) \times \mathbb{P}_1(C\langle\gamma\rangle)$  and is thus algebraic by Theorem 4.102. Since  $d' > d$ , Lemma 13.5 and Proposition 1.27 imply that  $(0: 1) \notin D_I$ . So  $D_I$  is a finite subset of  $\mathbb{P}_1(C\langle\gamma\rangle)$  by Lemma 4.101.

Thus the set of  $t \in C\langle\gamma\rangle$  such that  $k + 1$  polynomials each in  $P_i^*(1 - t, t)$ , have a common zero in  $C\langle\gamma\rangle^k$  is finite and its extension to  $C\langle\varepsilon, \delta, \gamma\rangle$  is a finite number of elements of  $C\langle\gamma\rangle$  which does not contain  $\delta$ .

It remains to prove that  $k$  polynomials  $Q_{i_j} \in \bar{P}_{i_j}^*$ ,  $j = 1, \dots, k$  have a finite number of common zeroes in  $R\langle\varepsilon, \delta, \gamma\rangle^k$ , which is an immediate consequence of Proposition 12.3, since  $d' > d$ .  $\square$

There is a close relationship between the sign conditions on  $\mathcal{P}$  and certain weak sign conditions on the polynomials in  $\mathcal{P}^*$  described by the following proposition. The role of the two infinitesimal quantities  $\delta$  and  $\gamma$  is the following:  $\delta$  is used to replace strict inequalities by weak inequalities and  $\gamma$  to replace equations by weak inequalities.



**Proposition 13.7.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{R}[X_1, \dots, X_k]$  be such that  $\deg P_i \leq d$  for all  $i$ , and suppose  $d' > d$ ,  $d'$  even. Let  $D \subset \mathbb{R}^k$  be a semi-algebraically connected component of the realization of the sign condition*

$$\begin{aligned} P_i &= 0, i \in I \subset \{1, \dots, s\}, \\ P_i &> 0, i \in \{1, \dots, s\} \setminus I. \end{aligned}$$

*Then there exists a semi-algebraically connected component  $D'$  of the subset  $\bar{D} \subset \mathbb{R}\langle \varepsilon, \delta, \gamma \rangle^k$  defined by the weak sign condition*

$$\begin{aligned} -\gamma \delta H_k(d', i) &\leq (1 - \delta) P_i \leq \gamma \delta H_k(d', i), \quad i \in I, \\ (1 - \delta) P_i &\geq \delta H_k(d', i), \quad i \in \{1, \dots, s\} \setminus I \\ \varepsilon^2 (X_1^2 + \dots + X_k^2) &\leq 1 \end{aligned}$$

*such that  $\lim_\gamma(D')$  is contained in the extension of  $D$  to  $\mathbb{R}\langle \varepsilon, \delta \rangle$ .*

**Proof:** If  $x \in D \subset \mathbb{R}^k$ , then  $x \in \bar{D}$ . Let  $D'$  be the semi-algebraically connected component of  $\bar{D}$  which contains  $x$ . Since  $\lim_\gamma$  is a ring homomorphism and  $d'$  is even, it is clear that  $\lim_\gamma(D')$  is contained in the realization of the conjunction of  $P_i = 0$ , for  $i \in I$ , and  $P_i > 0$ , for  $i \in \{1, \dots, s\} \setminus I$  in  $\mathbb{R}\langle \varepsilon, \delta \rangle^k$  and that it also contains  $x \in D$ . Since  $\bar{D}$  is bounded, by Proposition 12.43,  $\lim_\gamma(D')$  is also semi-algebraically connected. The statement of the proposition follows.  $\square$

**Corollary 13.8.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite subset of polynomials of degree less than  $d$  and suppose  $d' > d$ ,  $d'$  even. Let  $D$  be a semi-algebraically connected component of the realization of the sign condition*

$$\begin{aligned} P_i &= 0, i \in I \subset \{1, \dots, s\}, \\ P_i &> 0, i \in \{1, \dots, s\} \setminus I. \end{aligned}$$

*Then there exists a semi-algebraically connected component  $E'$  of the realization  $E \subset \mathbb{R}\langle \varepsilon, \delta, \gamma \rangle^{k+1}$  of*

$$\begin{aligned} -\gamma \delta H_k(d', i) &\leq (1 - \delta) P_i \leq \gamma \delta H_k(d', i), \quad 1i \in \{1, \dots, s\} \setminus I \in I, \\ (1 - \delta) P_i &\geq \delta H_k(d', i), \\ \varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) &= 1 \end{aligned}$$

*such that  $\Pi(\lim_\gamma(E'))$  is contained in the extension of  $D$  to  $\mathbb{R}\langle \varepsilon, \delta \rangle$ , where  $\Pi$  is the projection of  $\mathbb{R}^{k+1}$  to  $\mathbb{R}^k$  forgetting the last coordinate.*

As a consequence of Corollary 13.8, in order to compute all realizable sign conditions on  $\mathcal{P}$  it will be enough, using Proposition 13.1 and Proposition 13.6, to consider equations of the form

$$Q = Q_{i_1}^2 + \dots + Q_{i_j}^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2 = 0,$$

where  $j \leq k$ ,  $Q_{i_1} \in P_{i_1}^*, \dots, 1 \leq i_1 < \dots < i_j \leq s$ ,  $Q_{i_j} \in P_{i_j}^*$ , to find a point in each of the semi-algebraically connected components of their zero sets and to take their limit under  $\lim_\gamma$ .

A finite set  $\mathcal{S} \subset \mathbb{R}^k$  is a **set of sample points for  $\mathcal{P}$**  in  $\mathbb{R}^k$  if  $\mathcal{S}$  meets the realizations of all  $\sigma \in \text{SIGN}(\mathcal{P})$  (Notation 7.29). Note that the sample points output by Algorithm 11.2 (Cylindrical Decomposition) are a set of sample points for  $\mathcal{P}$  in  $\mathbb{R}^k$ , since the cells of a cylindrical decomposition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$  are  $\mathcal{P}$  invariant and partition  $\mathbb{R}^k$ . We are going to produce a set of sample points much smaller than the one output by Algorithm 11.2 (Cylindrical Decomposition), which was doubly exponential in the number of variables.

We present two versions of the Sampling algorithm. In the first one, the coordinates of the sample points belong to an extension of  $\mathbb{R}$  while in the second one the coordinates of the sample points belong to  $\mathbb{R}$ . The reason for presenting these two versions is technical: in Chapter 14, when we perform the same computation in a parametrized situation, the first version of Sampling will be easier to generalize, while in Chapter 15 it will be convenient to have sample points in  $\mathbb{R}^k$ . The two algorithms differ only in their last step.

We use the notation (12.18): let  $u = (f, g) \in \mathbb{K}[T]^{k+2}$ ,  $g = (g_0, \dots, g_k)$  be a  $k$ -univariate representation and  $P \in \mathbb{K}[X_1, \dots, X_k]$ .

$$P_u = g_0^e P\left(\frac{g_k}{g_0}, \dots, \frac{g_1}{g_0}\right), \tag{13.1}$$

where  $e$  is the least even number not less than the degree of  $P$ .

*Algorithm 13.1. [Computing Realizable Sign Conditions]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a set of  $s$  polynomials,

$$\mathcal{P} = \{P_1, \dots, P_s\} \subset D[X_1, \dots, X_k],$$

each of degree at most  $d$ .

- **Output:** a set of real univariate representations in  $D[\varepsilon, \delta, T]^{k+2}$  such that the associated points form a set of sample points for  $\mathcal{P}$  in  $\mathbb{R}\langle \varepsilon, \delta \rangle^k$ , meeting every semi-algebraically connected component of  $\text{Reali}(\sigma)$  for every  $\sigma \in \text{SIGN}(\mathcal{P})$  and the signs of the elements of  $\mathcal{P}$  at these points.
- **Complexity:**  $s^{k+1} d^{O(k)}$ .
- **Procedure:**
  - Initialize  $\mathcal{U}$  to the empty set.
  - Take as  $d'$  the smallest even natural number  $> d$ .
  - Define

$$\begin{aligned} P_i^* &= \{(1 - \delta) P_i + \delta H_k(d', i), (1 - \delta) P_i - \delta H_k(d', i), \\ &\quad (1 - \delta) P_i + \delta \gamma H_k(d', i), (1 - \delta) P_i - \delta \gamma H_k(d', i)\} \\ \mathcal{P}^* &= \{P_1^*, \dots, P_s^*\} \text{ for } 0 \leq i \leq s, \text{ using Notation 13.4.} \end{aligned}$$

- For every subset of  $j \leq k$  polynomials  $Q_{i_1} \in P_{i_1}^*, \dots, Q_{i_j} \in P_{i_j}^*$ ,
- Let

$$Q = Q_{i_1}^2 + \dots + Q_{i_j}^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2.$$

- For  $i = 1, \dots, k$ , let  $\bar{d}_i$  be the smallest even natural number greater than  $\deg_{X_i}(Q)$ ,  $i = 1, \dots, k$ , and let  $\bar{d}_{k+1} = 6$ ,  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k, \bar{d}_{k+1})$  and  $c = \varepsilon$ .
- Compute the multiplication table  $\mathcal{M}$  of  $\overline{\text{Cr}}(Q, \zeta)$  (Notation 12.46) using Algorithm 12.9 (Special Multiplication Table).
- Apply the  $\lim_{\gamma, \zeta}$  map using Algorithm 12.14 (Limit of Real Bounded Points) with input  $\mathcal{M}$ , and obtain a set of real univariate representations  $(v, \sigma)$  with

$$v = (f(T), g_0(T), \dots, g_k(T), g_{k+1}(T)) \in D[\varepsilon, \delta][T]^{k+3}.$$

- Ignore  $g_{k+1}(T)$  and consider only the real univariate representations  $(u, \sigma)$

$$u = (f(T), g_0(T), \dots, g_k(T)) \in D[\varepsilon, \delta][T]^{k+2}.$$

Add  $u$  to  $\mathcal{U}$ .

- Compute the signs of  $P \in \mathcal{P}$  at the points associated to the real univariate representations in  $\mathcal{U}$ , using Algorithm 10.13 (Univariate Sign Determination) with input  $f$  and its derivatives and the  $P_u$ ,  $P \in \mathcal{P}$ .

**Proof of correctness:** The correctness follows from Proposition 13.1, Proposition 12.42, Proposition 13.6, Corollary 13.8 and the correctness of Algorithm 12.9 (Special Multiplication Table), Algorithm 12.14 (Limit of Real Bounded Points) and Algorithm 10.13 (Univariate Sign Determination).  $\square$

**Complexity analysis:** The total number of  $j \leq k$ -tuples examined is  $\sum_{j \leq k} 4^j \binom{s}{j}$ . Hence, the number of calls to Algorithm 12.9 (Special Multiplication Table) and Algorithm 12.13 (Simple Univariate Representation) is also bounded by  $\sum_{j \leq k} 4^j \binom{s}{j}$ . Each such call costs  $d^{O(k)}$  arithmetic operations in  $D[\varepsilon, \delta, \gamma, \zeta]$ , using the complexity analysis of Algorithm 12.9 (Special Multiplication Table). Since there is a fixed number of infinitesimal quantities appearing with degree one in the input equations, the number of arithmetic operations in  $D$  is also  $d^{O(k)}$ , using the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table). Thus the total number of real univariate representations produced is bounded by  $\sum_{j \leq k} 4^j \binom{s}{j} O(d)^k$ , while the number of arithmetic operations performed for outputting sample points in  $R\langle \varepsilon, \delta \rangle^k$ , is bounded by  $\sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^k d^{O(k)}$ . The sign determination takes  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  arithmetic operations, using the complexity analysis of Algorithm 10.11 (Sign Determination).

If  $D = \mathbb{Z}$  and the bitsizes of the coefficients of the input polynomials are bounded by  $\tau$ , the size of the integer coefficients of the univariate representations are bounded by  $\tau d^{O(k)}$ , using the binary complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table).  $\square$

**Remark 13.9. [Hardness of the Existential Theory of the Reals]**

In computational complexity theory [127], the class of problems which can be solved in polynomial time is denoted by P. A problem is said to belong to the class NP if it can be solved by a non-deterministic Turing machine in polynomial time. Clearly,  $P \subset NP$  but it is unknown whether  $P \neq NP$ . A problem is NP-complete if it belongs to the class NP, and every other problem in NP can be reduced in polynomial time to an instance of this problem. A problem having only the latter property is said to be NP-hard. Since it is strongly believed that  $P \neq NP$ , it is very unlikely that an NP-hard problem will have a polynomial time algorithm.

It is a classical result in computational complexity that the Boolean satisfiability problem is NP-complete (see [127], Theorem 8.2 p. 171). The Boolean satisfiability problem is the following: given a Boolean formula,  $\phi(X_1, \dots, X_n)$ , written as a conjunction of disjunctions to decide whether it is satisfiable.

Since the Boolean satisfiability problem is NP-complete, it is very easy to see that the problem of existential theory of the reals is an NP-hard problem. Given an instance of a Boolean satisfiability problem, we can reduce it to an instance of the problem of the existential theory of the reals, by replacing each Boolean variable  $X_i$  by a real variable  $Y_i$  and adding the equation  $Y_i^2 - Y_i = 0$  and replacing each Boolean disjunction,  $X_{i_1} \vee X_{i_2} \vee \dots \vee X_{i_m}$  by the real inequality,

$$Y_{i_1} + Y_{i_2} + \dots + Y_{i_m} \geq 1.$$

It is clear that the original Boolean formula is satisfiable if and only if the semi-algebraic set defined by the corresponding real inequalities defined above is non-empty. This shows it is quite unlikely (unless  $P = NP$ ) that there exists any algorithm with binary complexity polynomial in the input size for the existential theory of the reals.  $\square$

**Remark 13.10. [Polynomial Space]**

Since Algorithm 13.1 (Computing Realizable Sign Conditions) is based essentially on computations of determinants of size  $O(d)^k$ , Remark 8.14 implies that it is possible to find the list of realizable sign conditions with complexity  $(s d)^{O(k)}$  using only  $(k \log_2(d))^{O(1)}$  amount of space at any time during the computation. In other words, the existential theory of the reals is in PSPACE [127]. The same remark applies to all the algorithms in Chapter 13, Chapter 15 and Chapter 16, as well as to all the algorithms in Chapter 14 when the number of block of quantifiers is bounded.  $\square$

*Algorithm 13.2.* [Sampling]

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a set of  $s$  polynomials,

$$\mathcal{P} = \{P_1, \dots, P_s\} \subset D[X_1, \dots, X_k],$$

each of degree at most  $d$ .

- **Output:** a set  $\mathcal{U}$  of real univariate representations in  $D[T]^{k+2}$  such that the associated points form a set of sample points for  $\mathcal{P}$  in  $R^k$ , meeting every semi-algebraically connected component of  $\text{Reali}(\sigma)$  for every  $\sigma \in \text{SIGN}(\mathcal{P})$ , and the signs of the elements of  $\mathcal{P}$  at these points.
- **Complexity:**  $s^{k+1} d^{O(k)}$ .
- **Procedure:**
  - Perform Algorithm 13.1 (Computing Realizable Sign Conditions) with input  $\mathcal{P}$  and output  $\mathcal{U}$ .
  - For every  $u \in \mathcal{U}$ ,  $u = (f, g)$ , replace  $\delta$  and  $\varepsilon$  by appropriately small elements from the field of quotients of  $D$  using Algorithm 11.20 (Removal of Infinitesimals) with input  $f$ , its derivatives and the  $P_u$ ,  $P \in \mathcal{P}$ . Then clear denominators to obtain univariate representation with entries in  $D[T]$ .

**Proof of correctness:** The correctness follows from the correctness of Algorithm 13.1 (Computing Realizable Sign Conditions) and Algorithm 11.20 (Removal of Infinitesimals)  $\square$

**Complexity analysis:** According to the complexity of Algorithm 13.1 (Computing Realizable Sign Conditions), the total number of real univariate representations produced is bounded by  $\sum_{j \leq k} 4^j \binom{s}{j} O(d)^k$ , while the number of arithmetic operations performed for outputting sample points in  $R\langle \varepsilon, \delta \rangle^k$ , is bounded by

$$\sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^k d^{O(k)}.$$

The sign determination takes  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  arithmetic operations.

Using Algorithm 11.20 (Removal of Infinitesimals) requires a further overhead of  $s d^{O(k)}$  arithmetic operations for every univariate representation output. Thus the number of arithmetic operations is bounded by

$$s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}.$$

However, the number of points actually constructed is only

$$\sum_{j \leq k} 4^j \binom{s}{j} O(d)^k.$$

If  $D = \mathbb{Z}$  and the bitsizes of the coefficients of the input polynomials are bounded by  $\tau$ , from the complexity of Algorithm 13.1 (Computing Realizable Sign Conditions), the size of the integer coefficients of the univariate representations in  $\mathcal{U}$  are bounded by  $\tau d^{O(k)}$ . In Algorithm 11.20 (Removal of Infinitesimals), we substitute a rational number, with numerator and denominator of bitsize  $\tau d^{O(k)}$ , in place of the variables  $\varepsilon$  and  $\delta$  and thus get points defined over  $\mathbb{Z}$  by polynomials with coefficients of bitsize  $\tau d^{O(k)}$ .  $\square$

Finally, we have proved the following theorem:

**Theorem 13.11.** *Let  $\mathcal{P}$  be a set of  $s$  polynomials each of degree at most  $d$  in  $k$  variables with coefficients in a real closed field  $R$ . Let  $D$  be the ring generated by the coefficients of  $\mathcal{P}$ . There is an algorithm that computes a set of  $\sum_{j \leq k} 4^j \binom{s}{j} O(d)^k$  points meeting every semi-algebraically connected component of the realization of every realizable sign condition on  $\mathcal{P}$  in  $R\langle \varepsilon, \delta \rangle^k$ . The algorithm has complexity  $\sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^k O(d)^k$  in  $D$ . There is also an algorithm computing the signs of all the polynomials in  $\mathcal{P}$  at each of these points with complexity  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  in  $D$ . The degrees of the univariate representations output are bounded by  $O(d)^k$ . If the polynomials in  $\mathcal{P}$  have coefficients in  $\mathbb{Z}$  of bitsizes at most  $\tau$ , the bitsizes of the coefficients of these univariate representations are bounded by  $\tau d^{O(k)}$ .*

Note that if we want the points to have coordinates in  $R^k$ , the complexity of finding sample points is also  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  in  $D$ , using Algorithm 13.2 (Sampling).

As a corollary,

**Theorem 13.12.** *Let  $R$  be a real closed field. Given a finite set,  $\mathcal{P} \subset R[X_1, \dots, X_k]$  of  $s$  polynomials each of degree at most  $d$ , then there exists an algorithm computing the set of realizable sign conditions  $\text{SIGN}(\mathcal{P})$  with complexity  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  in  $D$ , where  $D$  is the ring generated by the coefficients of the polynomials in  $\mathcal{P}$ .*

Recall that a  $\mathcal{P}$ -atom is one of  $P = 0, P \neq 0, P > 0, P < 0$ , where  $P$  is a polynomial in  $\mathcal{P}$  and a  $\mathcal{P}$ -formula is a formula written with  $\mathcal{P}$ -atoms. Since the truth or falsity of a sentence

$$(\exists X_1) \dots (\exists X_k) F(X_1, \dots, X_k),$$

where  $F(X_1, \dots, X_k)$  is a quantifier free  $\mathcal{P}$ -formula, can be decided by reading the list of realizable sign conditions on  $\mathcal{P}$ , the following theorem is an immediate corollary of Theorem 13.11.

**Theorem 13.13. [Existential Theory of the Reals]** *Let  $\mathbb{R}$  be a real closed field. Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite set of  $s$  polynomials each of degree at most  $d$ , and let*

$$(\exists X_1) \dots (\exists X_k) F(X_1, \dots, X_k),$$

*be a sentence, where  $F(X_1, \dots, X_k)$  is a quantifier free  $\mathcal{P}$ -formula. There exists an algorithm to decide the truth of the sentence with complexity  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$  in  $\mathbb{D}$  where  $\mathbb{D}$  is the ring generated by the coefficients of the polynomials in  $\mathcal{P}$ .*

*Remark 13.14.* Note that Theorem 13.11 implies that the total number of semi-algebraically connected components of realizable sign conditions defined by  $\mathcal{P}$  is bounded by  $\sum_{j \leq k} 4^j \binom{s}{j} O(d)^k$ . This bound is slightly less precise than the bound  $\sum_{j \leq k} \binom{s}{j} 4^j d (2d-1)^{k-1}$  given in Proposition 7.35, but does not require to use homology.  $\square$

## 13.2 A Few Applications

As a first application of the preceding results, we prove a bound on the radius of a ball meeting every semi-algebraically connected component of the realizations of the realizable sign conditions on a family of polynomials.

**Theorem 13.15.** *Given a set  $\mathcal{P}$  of  $s$  polynomials of degree at most  $d$  in  $k$  variables with coefficients in  $\mathbb{Z}$  of bitsizes at most  $\tau$ , there exists a ball of radius  $2^{\tau O(d)^k}$  intersecting every semi-algebraically connected component of the realization of every realizable sign condition on  $\mathcal{P}$ .*

**Proof:** The theorem follows from Theorem 13.11 together with Lemma 10.3 (Cauchy).  $\square$

We also have the following result.

**Theorem 13.16.** *Given a set  $\mathcal{P}$  of  $s$  polynomials of degree at most  $d$  in  $k$  variables with coefficients in  $\mathbb{Z}$  of bitsizes at most  $\tau$  such that  $S = \{x \in \mathbb{R}^k \mid P(x) > 0, P \in \mathcal{P}\}$  is a non-empty set, then in each semi-algebraically connected component of  $S$  there exists a point whose coordinates are rational numbers  $a_i/b_i$  where  $a_i$  and  $b_i$  have bitsizes  $\tau d^{O(k)}$ .*

**Proof:** This is a consequence of Theorem 13.11. We consider a point  $x$  belonging to  $S$  associated to a univariate representation

$$u = (f(T), g_0(T), \dots, g_k(T))$$

output by the algorithm, Algorithm 13.2 (Sampling), so that  $x_i = g_i(t)/g_0(t)$ , with  $t$  a root of  $f$  in  $\mathbb{R}$  known by its Thom encoding. Using Notation 13.8, each  $P_u(T)$  is of degree  $O(d)^k$ , and the bitsizes of its coefficients are bounded by  $\tau d^{O(k)}$ . Moreover, for every  $P \in \mathcal{P}$ ,  $P_u(t) > 0$ . Since the minimal distance between two roots of a univariate polynomial of degree  $O(d)^k$  with coefficients in  $\mathbb{Z}$  of bitsize  $\tau d^{O(k)}$  is at least  $2^{\tau d^{O(k)}}$  by Proposition 10.21, we get, considering polynomials  $P_u(T)$  for  $P \in \mathcal{P}$ , that there exists a rational number  $c/d$  with  $c$  and  $d$  of bitsizes  $\tau d^{O(k)}$  such that  $P_u(c/d)$  is positive. Thus defining the  $k$ -tuple  $a/b$  by  $a_i/b_i = (g_i/g_0)(c/d)$ , we get  $P(a/b) > 0$  for all  $P \in \mathcal{P}$  with bitsizes as claimed. □

We also apply our techniques to the algorithmic problem of checking whether an algebraic set has real dimension zero. We prove the following theorem.

Note that the only assumption we require in the second part of the theorem is that the real dimension of the algebraic set is 0 (the dimension of the complex part could be bigger).

**Theorem 13.17.** *Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$  have degree at most  $d$ , and let  $D$  be the ring generated by the coefficients of  $Q$ . There is an algorithm which checks if the real dimension of the algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  is 0 with complexity  $d^{O(k)}$  in  $D$ .*

*If the real dimension of  $\text{Zer}(Q, \mathbb{R}^k)$  is 0, the algorithm outputs a univariate representation of its points with complexity  $d^{O(k)}$  in  $D$ . Moreover, if  $D = \mathbb{Z}$  and the bitsizes of the coefficients are bounded by  $\tau$ , these points are contained in a ball of radius  $a/b$  with  $a$  and  $b$  in  $\mathbb{Z}$  of bitsizes  $\tau d^{O(k)}$ .*

*Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be  $s$  polynomials of degrees at most  $d$ , and let  $D$  the ring generated by the coefficients of the polynomials in  $\mathcal{P}$ . Then the signs of all the polynomials in  $\mathcal{P}$  at the points of  $\text{Zer}(Q, \mathbb{R}^k)$  can be computed with complexity  $s d^{O(k)}$  in  $D$ .*

**Proof:** In order to check whether the algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  is zero-dimensional, we apply Algorithm 13.1 (Sampling) to  $Q$ . Let  $\mathcal{U}$  be the set of real univariate representations output. Denote by  $K$  the finite set of points output, which intersects every semi-algebraically connected component of the algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$ . Now,  $\text{Zer}(Q, \mathbb{R}^k)$  is zero-dimensional, if and only if every point in  $K$  has a sufficiently small sphere centered around it, which does not intersect  $\text{Zer}(Q, \mathbb{R}^k)$ . For every  $(f(T), g_0(T), \dots, g_k(T), \sigma) \in \mathcal{U}$ , with associated point  $x$ , we introduce a new polynomial,

$$\begin{aligned}
 P(X_1, \dots, X_k, T) &= Q^2(X_1, \dots, X_k) + f^2(T) \\
 &\quad + ((g_0(T)X_1 - g_1(T))^2 \\
 &\quad + \dots \\
 &\quad + (g_0(T)X_k - g_k(T))^2 - g_0^2(T) \beta)^2
 \end{aligned}$$



where  $\beta$  is a new variable. We apply Algorithm 12.16 (Bounded Algebraic Sampling) to this  $(k + 1)$ -variate polynomial and check whether the corresponding zero set, intersected with the realization of the sign conditions  $\sigma$  on  $f(T)$  and its derivatives, is empty or not in  $\mathbb{R}\langle\beta\rangle^k$ .

If  $D = \mathbb{Z}$ , the bounds claimed follow from the fact that the polynomials in the univariate representations computed have integer coefficients of bit-sizes  $\tau d^{O(k)}$ .

For the third part we apply Algorithm 10.13 (Univariate Sign Determination) to the output of Algorithm 13.1 (Sampling). There are  $d^{O(k)}$  calls to Algorithm 10.13 (Univariate Sign Determination). The number of arithmetic operations is  $s d^{O(k)}$  in  $D$ .  $\square$

The following corollary follows immediately from the proof of Theorem 13.17.

**Corollary 13.18.** *Let  $\mathcal{Q}$  be a finite set of  $m$  polynomials in  $\mathbb{R}[X_1, \dots, X_k]$  of degree at most  $d$ . Then the coordinates of the isolated points of  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$  are zeros of polynomials of degrees  $O(d)^k$ . Moreover if  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials is bounded by  $\tau$ , then these points are contained in a ball of radius  $a/b$  with  $a$  and  $b$  in  $\mathbb{Z}$  of bitsizes  $(\tau + \log(m)) d^{O(k)}$  in  $\mathbb{R}^k$ .*

Our techniques can also be applied to decision problems in complex geometry.

**Proposition 13.19.** *Given a set  $\mathcal{P}$  of  $m$  polynomials of degree  $d$  in  $k$  variables with coefficients in  $\mathbb{C}$ , we can decide with complexity  $m d^{O(k)}$  in  $D$  (where  $D$  is the ring generated by the real and imaginary parts of the coefficients of the polynomials in  $\mathcal{P}$ ) whether  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  is empty. Moreover if  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials is bounded by  $\tau$ , then bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $(\tau + \log(m)) d^{O(k)}$ .*

**Proof:** Define  $Q$  as the sums of squares of the real and imaginary parts of the polynomials in  $\mathcal{P}$  and apply the Algorithm 13.1 (Computing realizable sign conditions).  $\square$

**Proposition 13.20.** *Given a set  $\mathcal{P}$  of  $m$  polynomials of degree  $d$  in  $k$  variables with coefficients in  $\mathbb{C}$ , we can decide with complexity  $m d^{O(k)}$  in  $D$  (where  $D$  is the ring generated by the real and imaginary parts of the coefficients of polynomials in  $\mathcal{P}$ ) whether the set of zeros of  $\mathcal{P}$  is zero-dimensional in  $\mathbb{C}^k$ . Moreover if  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then these points are contained in a ball of radius  $a/b$  with  $a$  and  $b$  in  $\mathbb{Z}$  of bitsize  $(\tau + \log(m)) d^{O(k)}$  in  $\mathbb{R}^k$ .*

**Proof:** Define  $Q$  as the sums of squares of the real and imaginary parts of the polynomials in  $\mathcal{P}$  and apply Theorem 13.17  $\square$

**Proposition 13.21.** *Given an algebraic set  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  where  $\mathcal{P}$  is a set of  $m$  polynomials of degree  $d$  in  $k$  variables with coefficients in  $\mathbb{C}$ , the real and imaginary parts of the coordinates of the isolated points of  $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$  are zeros of polynomials with coefficients in  $\mathbb{D}$  (where  $\mathbb{D}$  is the ring generated by the real and imaginary parts of the coefficients of polynomials in  $\mathcal{P}$ ) of degrees bounded by  $O(d)^k$ . Moreover if  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then these points are contained in a ball of radius  $a/b$  with  $a$  and  $b$  in  $\mathbb{Z}$  of bitsize  $(\tau + \log(m)) d^{O(k)}$  in  $\mathbb{R}^{2k} = \mathbb{C}^k$ .*

**Proof:** As in Proposition 13.18. □

### 13.3 Sample Points on an Algebraic Set

In this section we consider the problem of computing a set of sample points meeting the realizations of every realizable sign conditions of a family of polynomials restricted to an algebraic set. The goal is to have an algorithm whose complexity depends on the (possibly low) dimension of the algebraic set, which would be better than the complexity of Algorithms 13.1 and 13.2.

We prove the following theorem.

**Theorem 13.22.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be an algebraic set of real dimension  $k'$ , where  $Q$  is a polynomial in  $\mathbb{R}[X_1, \dots, X_k]$  of degree at most  $d$ , and let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite set of  $s$  polynomials with each  $P \in \mathcal{P}$  also of degree at most  $d$ . Let  $\mathbb{D}$  be the ring generated by the coefficients of  $Q$  and the polynomials in  $\mathcal{P}$ . There is an algorithm which computes a set of points meeting every semi-algebraically connected component of every realizable sign condition on  $\mathcal{P}$  in  $\text{Zer}(Q, \mathbb{R}_{\langle \varepsilon, \delta \rangle}^k)$ . The algorithm has complexity*

$$(k'(k - k') + 1) \sum_{j \leq k'} 4^j \binom{s}{j} d^{O(k)} = s^{k'} d^{O(k)}$$

*in  $\mathbb{D}$ . There is also an algorithm providing the signs of the elements of  $\mathcal{P}$  at these points with complexity*

$$(k'(k - k') + 1) \sum_{j \leq k'} 4^j \binom{s}{j} s d^{O(k)} = s^{k'+1} d^{O(k)}$$

*in  $\mathbb{D}$ . The degrees in  $T$  of the univariate representations output is  $O(d)^k$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$  then the bitsizes of the integers appearing in the intermediate computations and these univariate representations are bounded by  $\tau d^{O(k)}$ .*

*Remark 13.23.* This result is rather satisfactory since it fits with the bound on the number of realizable sign conditions proved in Proposition 7.35. □

We consider now a bounded algebraic set  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,  $c \leq 1$ , with  $Q \in \mathbb{R}[X_1, \dots, X_k]$ ,  $Q(x) \geq 0$  for every  $x \in \mathbb{R}^k$ , of degree bounded by  $d$  and of real dimension  $k'$ .

Let  $x$  be a smooth point of dimension  $k'$  of  $\text{Zer}(Q, \mathbb{R}^k)$ . We denote by  $T_x$  the tangent plane to  $\text{Zer}(Q, \mathbb{R}^k)$  at  $x$  and suppose that  $T_x$  is **transversal** to the  $k - k'$ -plane  $L$  defined by  $X_{k-k'+1} = \dots = X_k = 0$ , i.e. the intersection of  $T_x$  and  $L$  is  $\{0\}$ .

We will construct an algebraic set  $\text{Zer}(\mathcal{Q}, \mathbb{R}\langle\eta\rangle^k)$  covering  $\text{Zer}(Q, \mathbb{R}^k)$  in the neighborhood of  $x$ . The algebraic set  $\text{Zer}(\mathcal{Q}, \mathbb{R}\langle\eta\rangle^k)$  has two important properties. Firstly,  $\text{Zer}(\mathcal{Q}, \mathbb{R}\langle\eta\rangle^k)$  is defined by  $k - k'$  equations and has dimension  $k'$ . Secondly, any point  $z$  from a neighborhood of  $x$  in  $\text{Zer}(Q, \mathbb{R}^k)$  is infinitesimally close to a point from  $\text{Zer}(\mathcal{Q}, \mathbb{R}\langle\eta\rangle^k)$ .

Let  $\eta$  be a variable. Define, for  $d' > d$  and even, and  $\bar{d} = (d', \dots, d')$ ,

$$\begin{aligned} G(\bar{d}, c) &= c^{d'} X_1^{d'} + \dots + X_{k-k'}^{d'} + X_2^2 + \dots + X_{k-k'}^2 - (2(k - k') - 1) \\ \text{Def}(Q, \eta) &= \eta G(\bar{d}, c) + (1 - \eta) Q \\ \text{App}(Q, \eta) &= \left\{ \text{Def}(Q, \eta), \frac{\partial \text{Def}(Q, \eta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \eta)}{\partial X_{k-k'}} \right\}. \end{aligned}$$

**Lemma 13.24.**  $\dim(\text{Zer}(\text{App}(Q, \eta), \mathbb{R}\langle\eta\rangle^k)) \leq k'$ .

**Proof:** For every choice of  $z = (z_{k-k'+1}, \dots, z_k)$  in  $\mathbb{R}\langle\eta\rangle^{k'}$ , the affine  $(k - k')$ -plane  $L'$  defined by  $X_{k-k'+1} = z_{k-k'+1}, \dots, X_k = z_k$  intersects the algebraic set  $\text{Zer}(\text{App}(Q, \eta), \mathbb{C}^k\langle\eta\rangle)$  in at most a finite number of points. Indeed, consider the graded lexicographical ordering on the monomials for which  $X_1 < \dots < X_{k-k'}$ . Denoting  $\bar{X} = (X_1, \dots, X_{k-k'})$ , by Proposition 12.3,

$$\mathcal{G} = \text{App}(Q(\bar{X}, z), \eta)$$

is a Gröbner basis of the ideal  $\text{Ideal}(\mathcal{G}, \mathbb{R}\langle\eta\rangle)$  for the graded lexicographical ordering, since the leading monomials of elements of  $\mathcal{G}$  are pure powers of different  $X_i$ . Moreover, the quotient  $\mathbb{R}\langle\eta\rangle[\bar{X}]/\text{Ideal}(\mathcal{G}, \mathbb{R}\langle\eta\rangle)$  is a finite dimensional vector space and thus  $\mathcal{G}$  has a finite number of solutions in  $\mathbb{C}\langle\eta\rangle$  according to Proposition 12.7. The conclusion follows clearly by Corollary 5.28.  $\square$

**Proposition 13.25.** *There exists  $y \in \text{Zer}(\text{App}(Q, \eta), \mathbb{R}\langle\eta\rangle^k)$  such that  $\lim_{\eta} (y) = x$ .*

**Proof:** Since the tangent plane  $T_x$  to  $\text{Zer}(Q, \mathbb{R}^k)$  at  $x$  is transversal to  $L$ , the point  $x$  is an isolated point of the algebraic set

$$\text{Zer}(Q(\bar{X}, x_{k-k'+1}, \dots, x_k), \mathbb{R}\langle\eta\rangle^{k-k'}).$$

We can apply Proposition 12.42 to  $Q(\bar{X}, x_{k-k'+1}, \dots, x_k)$ .  $\square$

Using the preceding construction, we are able to approximate any smooth point such that  $T_x$  is transversal to the  $k - k'$ -plane defined by

$$X_{k-k'+1} = \dots = X_k = 0.$$

In order to approximate every point in  $\text{Zer}(Q, \mathbb{R}^k)$ , we are going to construct a family  $\mathcal{L}_{k, k-k'}$  of  $k - k'$ -planes with the following property: any linear subspace  $T$  of  $\mathbb{R}^k$  of dimension  $k'$  is transversal to at least one element of the family  $\mathcal{L}_{k, k-k'}$ , i.e. there is an element  $L$  of  $\mathcal{L}_{k, k-k'}$  such that  $T \cap L = \{0\}$ . The construction of  $\mathcal{L}_{k, k-k'}$  is based on properties of Vandermonde matrices.

**Notation 13.26.** Denoting by  $v_k(x)$  the Vandermonde vector  $(1, x, \dots, x^{k-1})$ , and by  $V_\ell$  the vector subspace of  $\mathbb{R}^k$  generated by  $v_k(\ell), v_k(\ell+1), \dots, v_k(\ell+k-k'-1)$ , it is clear that  $V_\ell$  is of dimension  $k - k'$  since the matrix of coordinates of vectors

$$v_{k-k'}(\ell), v_{k-k'}(\ell+1), \dots, v_{k-k'}(\ell+k-k'-1)$$

is an invertible Vandermonde matrix of dimension  $k - k'$ .

We now describe equations for  $V_\ell$ . Let, for  $k - k' + 1 \leq j \leq k$ ,

$$\begin{aligned} \bar{X}_j &= (X_1, \dots, X_{k-k'}, X_j), \\ v_{k-k', j}(\ell) &= (1, \dots, \ell^{k-k'-1}, \ell^{j-1}) \\ f_{\ell, j} &= \det(v_{k-k', j}(\ell), \dots, v_{k-k', j}(\ell+k-k'-1), \bar{X}_j), \\ L_{k', \ell}(X_1, \dots, X_k) &= (X_1, \dots, X_{k-k'}, f_{\ell, k-k'+1}, \dots, f_{\ell, k}). \end{aligned}$$

Note that the zero set of the linear forms  $f_{\ell, j}, k - k' + 1 \leq j \leq k$  is the vector space  $V_\ell$  and that  $L_{k', \ell}$  is a linear bijection such that  $L_{k', \ell}(V_\ell)$  consists of vectors of  $\mathbb{R}^k$  having their last  $k'$  coordinates equal to 0. We denote also by  $M_{k', \ell} = (d_{k-k', \ell})^{k'} L_{k', \ell}^{-1}$ , with

$$d_{k-k', \ell} = \det(v_{k-k'}(\ell), \dots, v_{k-k'}(\ell+k-k'-1)).$$

Note that  $M_{k', \ell}$  plays the same role as the inverse of  $L_{k', \ell}$  but is with integer coordinates, since, for  $k - k' + 1 \leq j \leq k$ ,  $d_{k-k', \ell}$  is the coefficient of  $X_j$  in  $f_{\ell, j}$ .

Define  $\mathcal{L}_{k, k-k'} = \{V_\ell \mid 0 \leq \ell \leq k'(k - k')\}$ .  $\square$

**Proposition 13.27.** Any linear subspace  $T$  of  $\mathbb{R}^k$  of dimension  $k'$  is transversal to at least one element of the family  $\mathcal{L}_{k, k-k'}$ .

**Corollary 13.28.** Any linear subspace  $T$  of  $\mathbb{R}^k$  of dimension  $j \geq k'$  is such that there exists  $0 \leq \ell \leq k'(k - k')$  such that  $V_\ell$  and  $T$  span  $\mathbb{R}^k$ .

In order to prove the proposition, we need the following lemma. Given a polynomial  $f(X) \in \mathbb{R}[X]$ , we denote by  $f^{\{n\}}(X)$  the  $n$ -th iterate of  $f$ , defined by

$$f^{\{0\}}(X) = X, f^{\{n+1\}}(X) = f(f^{\{n\}}(X)).$$

We denote by  $V^r(X)$  the vector subspace of  $\mathbb{R}(X)^k$  generated by

$$v_k(X), v_k(f(X)), \dots, v_k(f^{\{r-1\}}(X)).$$

By convention,  $V^0(X) = \{0\}$ .

**Lemma 13.29.** *Let  $T$  be a linear subspace of  $\mathbb{R}(X)^k$  of dimension  $\leq k'$ . Let  $f \in \mathbb{R}[X]$  be such that  $f^{\{i\}}(X) \neq f^{\{j\}}(X)$ , if  $i \neq j$ . Then the vector space  $V^{k-k'}(X)$  is transversal to  $T$  in  $\mathbb{R}(X)^k$ .*

**Proof:** The proof is by induction on  $k - k'$ . If  $k - k' = 0$ , the claim is clear since  $V^0(X) = \{0\}$ . Assume now by contradiction that  $k - k' \geq 1$  and  $V^{k-k'}(X)$  is not transversal to  $T$ . By induction hypothesis,  $V^{k-k'-1}(X)$  is transversal to  $T$ . Hence  $v(f^{\{k-k'-1\}}(X))$  belongs to the vector space generated by  $T$  and  $V^{k-k'-1}(X)$ .

It follows by induction on  $j$  that for every  $j \geq k - k'$ ,  $v(f^{\{j-1\}}(X))$  belongs to the vector space generated by  $T$  and  $V^{k-k'-1}(X)$ . Consider the Vandermonde matrix  $V(X, \dots, f^{\{k-1\}}(X))$ . Since

$$\det(V(X, \dots, f^{\{k-1\}}(X))) = \prod_{k-1 \geq i > j \geq 0} (f^{\{i\}}(X) - f^{\{j\}}(X)) \neq 0,$$

and the dimension of the vector space generated by  $T$  and  $V^{k-k'-1}(X)$  is  $< k$ , we obtained a contradiction.  $\square$

**Proof of Proposition 13.27:** We apply Lemma 13.29 to  $f(X) = X + 1$ . Denoting by  $e_1, \dots, e_{k'}$  a basis of  $T$ , and applying Lemma 13.29

$$D = \det(e_1, \dots, e_{k'}, v_k(X), v_k(X+1), \dots, v_k(X+k-k'-1))$$

is not identically 0. Since

$$\begin{aligned} D' &= \det(e_1, \dots, e_{k'}, v_k(X_1), v_k(X_2), \dots, v_k(X_{k-k'})) \\ &= \left( \prod_{1 \leq i < j \leq k-k'} (X_i - X_j) \right) S(X_1, \dots, X_{k-k'}) \end{aligned}$$

with  $S(X_1, \dots, X_{k-k'}) \in \mathbb{R}[X_1, \dots, X_{k-k'}]$ , and the degree of  $D'$  is bounded by  $\sum_{k'}^{k-1} i = (1/2)(k - k')(k + k' - 1)$  the degree of  $S$  is bounded by  $(1/2)(k - k')(k + k' - 1) - \binom{k-k'}{2} = k'(k - k')$ . Since  $(X + i) - (X + j)$  is a constant, it is clear that the degree of  $D = S(X, \dots, X + k - k' - 1)$  is also bounded by  $k'(k - k')$ . Hence, there exists  $\ell \in \{0, \dots, k'(k - k')\}$  which is not a root of  $D$ . The corresponding  $V_\ell$  is transversal to  $T$ .  $\square$

**Notation 13.30.** Let  $\eta$  be a variable,  $d'$  an even natural number,  $d' > d$ , and  $\bar{d} = (d', \dots, d')$  and  $0 \leq \ell \leq k'(k - k')$ . Using Notation 13.26, let

$$\begin{aligned} Q_\ell(X_1, \dots, X_k) &= Q(M_{k', \ell}(X_1, \dots, X_k)), \\ \text{Def}(Q_\ell, \eta) &= \eta G(\bar{d}, c) + (1 - \eta) Q_\ell, \\ \text{App}(Q_\ell, \eta) &= \left\{ \text{Def}(Q_\ell, \eta), \frac{\partial \text{Def}(Q_\ell, \eta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q_\ell, \eta)}{\partial X_{k-k'}} \right\}. \end{aligned}$$

Define  $Z_\ell = M_{k',\ell}(\text{Zer}(\text{App}(Q_\ell, \eta), \mathbb{R}\langle \eta \rangle^k))$ ,  $0 \leq \ell \leq k'(k - k')$ , the **approximating varieties** of  $\text{Zer}(Q, \mathbb{R}^k)$ .  $\square$

This terminology is justified by the following result.

**Proposition 13.31.**

$$\lim_{\eta} \left( \bigcup_{\ell=0}^{k'(k-k')} Z_\ell \right) = \text{Zer}(Q, \mathbb{R}^k).$$

**Proof:** It is clear that

$$\lim_{\eta} (\text{Zer}(\text{App}(Q_\ell, \eta), \mathbb{R}\langle \eta \rangle^k)) \subset \text{Zer}(Q_\ell, \mathbb{R}^k).$$

Denote by  $S_\ell$  the set of all smooth points of  $\text{Zer}(Q, \mathbb{R}^k)$  having a tangent  $k''$ -plane ( $k'' \leq k'$ ) transversal to  $V_\ell$ . Using Notation 13.26, Proposition 13.25 and Proposition 13.27, the union of the  $S_\ell$  for  $0 \leq \ell \leq k'(k - k')$  is a semi-algebraic subset of  $\text{Zer}(Q, \mathbb{R}^k)$  whose closure is  $\text{Zer}(Q, \mathbb{R}^k)$ , using Proposition 5.54. The image under  $\lim_{\eta}$  of  $\bigcup_{\ell=0}^{k'(k-k')} Z_\ell$  contains  $\bigcup_{\ell=0}^{k'(k-k')} S_\ell$ . Moreover, by Proposition 12.43, the image under  $\lim_{\eta}$  of a semi-algebraic set is closed. Hence,

$$\lim_{\eta} \left( \bigcup_{\ell=0}^{k'(k-k')} Z_\ell \right) \supset \text{Zer}(Q, \mathbb{R}^k). \quad \square$$

**Notation 13.32.** Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset \mathbb{R}[X_1, \dots, X_k]$  be polynomials of degree bounded by  $d$ ,  $d'$  an even natural number,  $d' > d$ ,  $\bar{d} = (d', \dots, d')$ . In order to perturb the polynomials in  $\mathcal{P} = \{P_1, \dots, P_s\}$  to get a family in  $k'$ -general position with  $\text{App}(Q, \eta)$ , we use polynomials  $H(d'', i) = 1 + \sum_{1 \leq j \leq k-k'} i^j X_{k-k'+j}^{d''}$ , for  $1 \leq i \leq s$ , where  $d''$  is an even natural number,  $d'' > d'$ . We consider two variables  $\delta, \gamma$ . Let

$$\begin{aligned} P_i^* &= \{(1 - \delta) P_i + \delta H(d'', i), (1 - \delta) P_i - \delta H(d'', i), \\ &\quad (1 - \delta) P_i + \gamma \delta H(d'', i), (1 - \delta) P_i - \gamma \delta H(d'', i)\} \\ \mathcal{P}^* &= \{P_1^*, \dots, P_s^*\}. \end{aligned}$$

$\square$

**Proposition 13.33.** *The family  $\mathcal{P}^*$  is in strong  $k'$ -general position with respect to  $\text{App}(Q, \eta)$  in  $\mathbb{R}\langle \delta, \gamma, \eta \rangle^k$ .*

The proof of the proposition uses the following lemma.

**Lemma 13.34.** *The polynomials  $H(d'', i)$ ,  $0 \leq i \leq s$ , are in  $k'$ -general position with respect to  $\text{App}(Q, \eta)$  in  $\mathbb{R}\langle \eta \rangle^k$ .*

**Proof:** Let

$$\begin{aligned} \text{Def}(Q, \lambda, \mu) &= \lambda Q + \mu G(\bar{d}, c) \\ \text{App}(Q, \mu) &= \left\{ \text{Def}(Q, \lambda, \mu), \frac{\partial \text{Def}(Q, \lambda, \mu)}{\partial X_{k'+2}}, \dots, \frac{\partial \text{Def}(Q, -\lambda, \mu)}{\partial X_k} \right\} \\ \text{App}^h(Q, \lambda, \mu) &= \left\{ \text{Def}^h(Q, \lambda, \mu), \frac{\partial \text{Def}^h(Q, \lambda, \mu)}{\partial X_{k'+2}}, \dots, \frac{\partial \text{Def}^h(Q, \lambda, \mu)}{\partial X_k} \right\} \end{aligned}$$

where

$$\text{Def}^h(Q, \lambda, \mu) = X_0^{d'} \text{Def}(Q, \lambda, \mu) \left( \frac{X_1}{X_0}, \dots, \frac{X_k}{X_0} \right).$$

The system  $\text{App}^h(Q, 0, 1)$  has only the solution  $(1: 0: \dots: 0)$  in  $\mathbb{P}_{k-k'}(\mathbb{C})$ . The polynomials  $H^h(d'', i)$ ,  $0 \leq i \leq s$ , are in  $k'$ -general position in  $\mathbb{P}_{k'}(\mathbb{C})$  by Lemma 13.5. Thus, with  $J = \{j_1, \dots, j_{k'+1}\} \subset \{1, \dots, s\}$ , the set  $D_J$  of  $(\lambda: \mu) \in \mathbb{P}_1(\mathbb{C})$  such that  $H^h(d'', j_1), \dots, H^h(d'', j_{k'+1})$  have a common zero on  $\text{App}^h(Q, \lambda, \mu)$  does not contain  $(0: 1)$ .

Moreover, the set  $D_J$  is the projection to  $\mathbb{P}_1(\mathbb{C})$  of an algebraic subset of  $\mathbb{P}_k(\mathbb{C}) \times \mathbb{P}_1(\mathbb{C})$  and is thus algebraic by Theorem 4.102. Since  $d'' > d'$  and  $(0: 1) \notin D_J$ ,  $D_J$  is a finite subset of  $\mathbb{P}_1(\mathbb{C})$ . Thus the set of  $t \in \mathbb{C}$  such that  $k' + 1$  polynomials among  $H_T(d'', j)$ ,  $j \leq s$ , have a common zero on  $\text{Zer}(\text{App}(Q, t, 1 - t), \mathbb{C}^k)$  is finite, and its extension to  $\mathbb{C}\langle \eta \rangle$  is a finite set of elements which does not contain  $\eta$ .  $\square$

**Proof of Proposition 13.33:** Consider

$$\begin{aligned} \bar{P}_i^{*h} &= \{ \lambda P_i^h + \mu H^h(d'', i), \lambda P_i^h - \mu H^h(d'', i), \\ &\quad \lambda P_i^h + \mu \gamma H^h(d'', i), \lambda P_i^h - \mu \gamma H^h(d'', i) \}, \end{aligned}$$

$0 \leq i \leq s$ . Let  $J = \{j_1, \dots, j_{k'+1}\} \subset \{1, \dots, s\}$  and  $A_{j_i} \in \bar{P}_i^{*h}$ . The set  $D_J$  of  $(\lambda: \mu)$  such that  $A_{j_1}(\lambda, \mu), \dots, A_{j_{k'+1}}(\lambda, \mu)$  have a common zero with

$$\text{Zer}(\text{App}^h(Q, \lambda, \mu), \mathbb{P}_k(\mathbb{C}))$$

in  $\mathbb{P}_k(\mathbb{C}\langle \gamma, \eta \rangle)$  is the projection to  $\mathbb{P}_1(\mathbb{C}\langle \gamma, \eta \rangle)$  of an algebraic subset of  $\mathbb{P}_k(\mathbb{C}\langle \gamma, \eta \rangle) \times \mathbb{P}_1(\mathbb{C}\langle \gamma, \eta \rangle)$  and is thus algebraic by Theorem 4.102. Since  $d'' > d'$ ,  $(0: 1) \notin D_J$  by Lemma 13.34, and  $D_J$  is a finite subset of  $\mathbb{P}_1(\mathbb{C}\langle \gamma, \eta \rangle)$ . Thus the set of  $t \in \mathbb{C}\langle \gamma, \eta \rangle$  such that  $k' + 1$  polynomials among  $(1 - t)P_i + tH(d'', j)$ ,  $j \leq s$ , have a common zero on

$$\text{Zer}(\text{App}(Q, \eta), \mathbb{R}\langle \gamma, \eta \rangle^k)$$

is finite, and its extension to  $\mathbb{C}\langle \delta, \gamma, \eta \rangle$  is a finite set of elements of  $\mathbb{C}\langle \delta, \gamma, \eta \rangle$  which does not contain  $\delta$ . It remains to prove that  $k'$  polynomials

$$A_{j_1}(\lambda, \mu), \dots, A_{j_{k'}}(\lambda, \mu)$$

have a finite number of common zeroes in  $\mathbb{R}\langle \delta, \eta, \gamma \rangle^k$ , which is an immediate consequence of Proposition 12.3, since  $d' > d$ .  $\square$

We consider now a polynomial  $Q \in \mathbb{R}[X_1, \dots, X_k]$ , with  $\text{Zer}(Q, \mathbb{R}^k)$  not necessarily bounded.

The following proposition holds.

**Proposition 13.35.** *Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$  and  $\mathcal{P} = \{P_1, \dots, P_s\}$  be a finite subset of  $\mathbb{R}[X_1, \dots, X_k]$ . Let  $d$  be a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ ,  $d'$  an even number  $> 2d$ , and  $d''$  an even number  $> d'$ . Let  $D$  be a connected component of the realization of the sign condition*

$$\begin{aligned} Q &= 0 \\ P_i &= 0, i \in I \subset \{1, \dots, s\} \\ P_i &> 0, i \in \{1, \dots, s\} \setminus I. \end{aligned}$$

Let  $\bar{Q} = Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2$ . If the set  $E \subset \mathbb{R}(\varepsilon, \delta, \gamma, \eta)^k$  described by

$$\begin{aligned} \bigwedge_{R \in \text{App}(\bar{Q}, \eta, \bar{d}, c)} R &= 0 \\ -\gamma \delta H_k(\bar{d}, i) &\leq (1 - \delta) P_i \leq \gamma \delta H_k(\bar{d}, i), i \in I, \\ (1 - \delta) P_i &\geq \delta H_k(\bar{d}, i), i \in \{1, \dots, s\} \setminus I \\ \varepsilon^2 (X_1^2 + \dots + X_k^2) &\leq 1 \end{aligned}$$

is non-empty, there exists a connected component  $E'$  of  $E$  such that  $\pi(\lim_{\gamma, \eta} (E'))$  is contained in the extension of  $D$  to  $\mathbb{R}(\varepsilon, \delta)$ , where  $\pi$  is the projection of  $\mathbb{R}^{k+1}$  to  $\mathbb{R}^k$  forgetting the last coordinate.

**Proof:** The proof is similar to the proof of Proposition 13.7, using Proposition 13.31. □

**Notation 13.36.** The set  $\text{SIGN}(\mathcal{P}, \mathcal{Q}) \subset \{0, 1, -1\}^{\mathcal{P}}$  is the set of all realizable sign conditions for  $\mathcal{P}$  on  $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ . If  $\sigma \in \text{SIGN}(\mathcal{P}, \mathcal{Q})$  we denote

$$\text{Reali}(\sigma, \mathcal{Q}) = \{x \in \mathbb{R}^k \mid \bigwedge_{Q \in \mathcal{Q}} Q(x) = 0 \wedge \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma\}.$$

For  $0 \leq \ell \leq k'(k - k')$ , and  $P \in \mathbb{R}[X_1, \dots, X_k]$  we denote by

$$P_\ell(X_1, \dots, X_k) = P(M_{k', \ell}(X_1, \dots, X_k))$$

If  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ ,  $\mathcal{P}_\ell = \{P_\ell \mid P \in \mathcal{P}\}$ .

Given a real univariate representation

$$v = (f(T), g_0(T), g_1(T), \dots, g_k(T)), \sigma),$$

with associated point  $z$ , we denote by

$$M_{k', \ell}(v) = (f(T), g_0(T), h_1(T), \dots, h_k(T)),$$



with  $h_1(T), \dots, h_k(T) = M_{k',\ell}(g_1(T), \dots, g_k(T))$ , the real univariate representation with associated point  $M_{k',\ell}(z)$ .  $\square$

*Algorithm 13.3. [Sampling on an Algebraic Set]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  of degree at most  $d$ , with  $\text{Zer}(Q, R^k)$  of real dimension  $k'$ ,
  - a set of  $s$  polynomials,  $\mathcal{P} = \{P_1, \dots, P_s\} \subset D[X_1, \dots, X_k]$ , each of degree at most  $d$ .
- **Output:** a set  $\mathcal{U}$  of real univariate representations in  $D[\varepsilon, \delta][T]^{k+2}$  such that for every  $\sigma \in \text{SIGN}(\mathcal{P}, Q)$ , the associated points meet every semi-algebraically connected component of the extension of  $\text{Reali}(\sigma, Q)$  to  $R\langle\varepsilon, \delta\rangle^k$ .
- **Complexity:**  $s^{k'+1}d^{O(k)}$ .
- **Procedure:**
  - Take  $d' = 2(d+1), \bar{d} = (d', \dots, d'), d'' = 2(d+2)$ .
  - For every  $0 \leq \ell \leq k'(k-k')$ , define

$$\bar{Q}_\ell = Q_\ell^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2,$$

and define  $\text{App}(\bar{Q}_\ell, \eta, \bar{d}, \varepsilon)$  and  $\mathcal{P}_\ell^*$ , using Notation 13.30 and Notation 13.32.

- For every  $j \leq k'$ -tuple of polynomials  $A_{t_1} \in P_{t_1}^*, \dots, A_{t_j} \in P_{t_j}^*$  let

$$R = \sum_{P \in \text{App}(\bar{Q}_\ell, \eta, \bar{d}, \varepsilon)} P^2 + A_{i_1}^2 + \dots + A_{i_j}^2.$$

- Take for  $i = 1, \dots, k$ ,  $\bar{d}_i$  equal to the smallest even natural number  $> \deg_{X_i}(R)$ ,  $\bar{d}_{k+1} = 8$ ,  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k, \bar{d}_{k+1})$ ,  $c = \varepsilon$ .
- Compute the multiplication table  $\mathcal{M}$  of  $\overline{\text{Cr}}(R, \zeta)$  (Notation 12.46) using Algorithm 12.9 (Special Multiplication Table). Apply the  $\lim_{\gamma, \eta, \zeta}$  map using Algorithm 12.14 (Limit of Real Bounded Points) with input  $\mathcal{M}$ , and obtain a set  $\mathcal{U}_\ell$  of real univariate representations  $v$  with

$$v = ((f(T), g_0(T), \dots, g_k(T)), \sigma) \\ (f(T), g_0(T), \dots, g_k(T)) \in D[\varepsilon, \delta][T]^{k+2}.$$

- Define  $\mathcal{U} = \bigcup_{\ell=0}^{k'(k-k')} M_{k',\ell}(\mathcal{U}_\ell)$ . Compute the signs of  $P \in \mathcal{P}$  at the points associated to the real univariate representations  $v$  in  $\mathcal{U}$ ,

$$v = (f(T), g_0(T), \dots, g_k(T)), \sigma)$$

using Algorithm 10.13 (Univariate Sign Determination) with input  $f$  and its derivatives and  $\mathcal{P}$ .

**Proof of correctness:** Follows from Proposition 13.1, Proposition 13.31, Proposition 13.33, and Proposition 13.35. □

**Complexity analysis:** It is clear that  $\sum_{j \leq k'} 4^j \binom{s}{j}$  tuples of polynomials are considered for each  $0 \leq \ell \leq k'(k - k')$ . The cost for each such tuple is  $d^{O(k)}$  using the complexity analysis of Algorithm 12.18 (Parametrized Bounded Algebraic Sampling), since we are using a fixed number of infinitesimal quantities. Hence, the complexity for finding sample points in  $\mathbb{R}\langle \varepsilon, \delta \rangle$  is bounded by  $(k'(k - k') + 1) \sum_{j \leq k'} 4^j \binom{s}{j} d^{O(k)} = s^{k'} d^{O(k)}$ . Note that the degrees of the polynomials output are bounded by  $O(d)^k$  and that when  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of  $Q$  and  $P \in \mathcal{P}$  are bounded by  $\tau$ , the bitsizes of the coefficients of the polynomials occurring in the multiplication table are  $\tau d^{O(k)}$ . Moreover the number of real univariate representations output is  $s^{k'} O(d)^k$ .

The cost of computing the signs is  $s d^{O(k)}$  per point associated to a real univariate representation. Hence, the complexity of the sign determination at the end of the algorithm is bounded by

$$(k'(k - k') + 1) \sum_{j \leq k'} 4^j \binom{s}{j} s d^{O(k)} = s^{k'+1} d^{O(k)}.$$

Note that if we want the points to have coordinates in  $\mathbb{R}^k$ , the complexity of finding sample points is still  $s^{k'+1} d^{O(k)}$  in  $D$ , using Algorithm 11.20 (Removal of Infinitesimals).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ . □

**Proof of Theorem 13.22:** The claim is an immediate consequence of the complexity analysis of Algorithm 13.3 (Sampling on an Algebraic Set). □

*Remark 13.37.* The complexity of Algorithm 13.3 is rather satisfactory since it fits with the bound on the number of realizable sign conditions proved in Proposition 7.35. □

The following result is an immediate corollary of Theorem 13.22.

**Theorem 13.38.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be an algebraic set of real dimension  $k'$ , where  $Q$  is a polynomial in  $\mathbb{R}[X_1, \dots, X_k]$  of degree at most  $d$ , and let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite set of  $s$  polynomials with each  $P \in \mathcal{P}$  also of degree at most  $d$ . Let  $D$  be the ring generated by the coefficients of  $Q$  and the polynomials in  $\mathcal{P}$ . There is an algorithm that takes as input  $Q, k',$  and  $\mathcal{P}$  and computes  $\text{SIGN}(\mathcal{P}, Q)$  with complexity*

$$(k'(k - k') + 1) \sum_{j \leq k'} 4^j \binom{s}{j} s d^{O(k)} = s^{k'+1} d^{O(k)}$$

in  $D$ . If  $D=Z$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ .

*Remark 13.39.* Note that the dimension of the algebraic set is part of the input. A method for computing the dimension of an algebraic set is given at the end of Chapter 14.  $\square$

## 13.4 Computing the Euler-Poincaré Characteristic of Sign Conditions

Our aim is to give a method for determining the Euler-Poincaré characteristic of the realization of sign conditions realized by a finite set  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  on an algebraic set  $Z = \text{Zer}(Q, \mathbb{R}^k)$ , with  $Q \in \mathbb{R}[X_1, \dots, X_k]$ .

This is done by a method very similar to Algorithm 10.11 (Sign Determination): we compute Euler-Poincaré characteristics of realizations of sign conditions rather than cardinalities of sign conditions on a finite set, using the notion of Euler-Poincaré-query rather than that of Tarski-query.

We recall the following definitions already introduced in Section 6.3.

Given  $S$  a locally closed semi-algebraic set contained in  $Z$ , we denote by  $\chi(S)$  the Euler-Poincaré characteristic of  $S$ .

Given  $P \in \mathbb{R}[X_1, \dots, X_k]$ , we denote

$$\begin{aligned} \text{Reali}(P=0, S) &= \{x \in S \mid P(x) = 0\}, \\ \text{Reali}(P > 0, S) &= \{x \in S \mid P(x) > 0\}, \\ \text{Reali}(P < 0, S) &= \{x \in S \mid P(x) < 0\}, \end{aligned}$$

and  $\chi(P=0, S)$ ,  $\chi(P > 0, S)$ ,  $\chi(P < 0, S)$  the Euler-Poincaré characteristics of the corresponding sets. The Euler-Poincaré-query of  $P$  for  $S$  is

$$\text{EuQ}(P, S) = \chi(P > 0, S) - \chi(P < 0, S).$$

Let  $\mathcal{P} = P_1, \dots, P_s$  be a finite list of polynomials in  $\mathbb{R}[X_1, \dots, X_k]$ .

Let  $\sigma$  be a sign condition on  $\mathcal{P}$ . The **realization of the sign condition  $\sigma$  at  $S$**  is

$$\text{Reali}(\sigma, S) = \{x \in S \mid \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma(P)\},$$

and its Euler-Poincaré characteristic is denoted  $\chi(\sigma, S)$ .

**Notation 13.40.** Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$ ,  $Z = \text{Zer}(Q, \mathbb{R}^k)$ . We denote as usual by  $\text{SIGN}(\mathcal{P}, Z)$  the list of  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$  such that  $\text{Reali}(\sigma, Z)$  is non-empty. We denote by  $\chi(\mathcal{P}, Z)$  the list of Euler-Poincaré characteristics  $\chi(\sigma, Z) = \chi(\text{Reali}(\sigma, Z))$  for  $\sigma \in \text{SIGN}(\mathcal{P}, Z)$ . We are going to compute  $\chi(\mathcal{P}, Z)$ , using Euler-Poincaré-queries of products of elements of  $\mathcal{P}$ .  $\square$

We use Notation 10.67, and order lexicographically  $\{0, 1, -1\}^{\mathcal{P}}$  and  $\{0, 1, 2\}^{\mathcal{P}}$ . Given  $A = \alpha_1, \dots, \alpha_m$  a list of elements of  $\{0, 1, 2\}^{\mathcal{P}}$ , with  $\alpha_1 <_{\text{lex}} \dots <_{\text{lex}} \alpha_m$ , we write  $\mathcal{P}^A$  for  $\mathcal{P}^{\alpha_1}, \dots, \mathcal{P}^{\alpha_m}$ , and  $\text{EuQ}(\mathcal{P}^A, S)$  for  $\text{EuQ}(\mathcal{P}^{\alpha_1}, S), \dots, \text{EuQ}(\mathcal{P}^{\alpha_m}, S)$ .

We denote by  $\text{Mat}(A, \Sigma)$  the matrix of signs of  $\mathcal{P}^A$  on  $\Sigma$  (see Definition 10.3).

**Proposition 13.41.** *If  $\cup_{\sigma \in \Sigma} \text{Reali}(\sigma, S) = S$ , then*

$$\text{Mat}(A, \Sigma) \cdot \chi(\Sigma, S) = \text{EuQ}(\mathcal{P}^A, S).$$

**Proof:** The proof is by induction on the number  $s$  of polynomials in  $\mathcal{P}$ . The statement when  $s = 1$  follows from Proposition 6.60, since the Euler-Poincaré characteristic of an empty sign condition is zero. Suppose the statement holds for  $\mathcal{P}' = P_1, \dots, P_{s-1}$  and consider  $\mathcal{P} = P_1, \dots, P_s$ . Define

$$\begin{aligned} \Sigma_0 &= \{\sigma \in \Sigma \mid \sigma(P_s) = 0\} \\ \Sigma_1 &= \{\sigma \in \Sigma \mid \sigma(P_s) = 1\} \\ \Sigma_{-1} &= \{\sigma \in \Sigma \mid \sigma(P_s) = -1\}, \\ T &= \bigcup_{\sigma \in \Sigma_0} \text{Reali}(\sigma, S) \\ U &= \bigcup_{\sigma \in \Sigma_1} \text{Reali}(\sigma, S) \\ V &= \bigcup_{\sigma \in \Sigma_{-1}} \text{Reali}(\sigma, S). \end{aligned}$$

Note that  $T, U$ , and  $V$  are all locally closed whenever  $S$  is locally closed. Let  $\alpha \in \{0, 1, 2\}^{\mathcal{P}}$  and  $\alpha' \in \{0, 1, 2\}^{\mathcal{P}'}$  defined by  $\alpha'(P_j) = \alpha(P_j)$ , for  $1 \leq j \leq s-1$ . Using the additive property of Euler-Poincaré characteristic (Proposition 6.57),

$$\begin{aligned} \chi(\mathcal{P}^\alpha = 0, S) &= \chi(\mathcal{P}^\alpha = 0, T) + \chi(\mathcal{P}^\alpha = 0, U) + \chi(\mathcal{P}^\alpha = 0, V), \\ \chi(\mathcal{P}^\alpha > 0, S) &= \chi(\mathcal{P}^\alpha > 0, T) + \chi(\mathcal{P}^\alpha > 0, U) + \chi(\mathcal{P}^\alpha > 0, V), \\ \chi(\mathcal{P}^\alpha < 0, S) &= \chi(\mathcal{P}^\alpha < 0, T) + \chi(\mathcal{P}^\alpha < 0, U) + \chi(\mathcal{P}^\alpha < 0, V). \end{aligned}$$

– If  $\alpha(P_s) = 0$ ,

$$\text{EuQ}(\mathcal{P}^\alpha, S) = \text{EuQ}(\mathcal{P}'^{\alpha'}, T) + \text{EuQ}(\mathcal{P}'^{\alpha'}, U) + \text{EuQ}(\mathcal{P}'^{\alpha'}, V).$$

– If  $\alpha(P_s) = 1$ ,

$$\text{EuQ}(\mathcal{P}^\alpha, S) = \text{EuQ}(\mathcal{P}'^{\alpha'}, U) - \text{EuQ}(\mathcal{P}'^{\alpha'}, V).$$

– If  $\alpha(P_s) = 2$ ,

$$\text{EuQ}(\mathcal{P}^\alpha, S) = \text{EuQ}(\mathcal{P}'^{\alpha'}, U) + \text{EuQ}(\mathcal{P}'^{\alpha'}, V).$$

The claim follows from the induction hypothesis applied to  $T, U$  and  $V$ , the definition of  $\text{Mat}(A, \Sigma)$  (Definition 2.66) and the additive property of Euler-Poincaré characteristic (Proposition 6.57), which implies, for every  $\sigma \in \Sigma$ ,

$$\chi(\sigma, S) = \chi(\sigma, T) + \chi(\sigma, U) + \chi(\sigma, V). \quad \square$$

Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$ ,  $Z = \text{Zer}(Q, \mathbb{R}^k)$ . We consider a list  $A(Z)$  of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  **adapted to sign determination** for  $\mathcal{P}$  on  $Z$ , i.e. such that the matrix of signs of  $\mathcal{P}^A$  over  $\text{SIGN}(\mathcal{P}, Z)$  is invertible. If  $\mathcal{P} = P_1, \dots, P_s$ , let  $\mathcal{P}_i = P_i, \dots, P_s$ , for  $0 \leq i \leq s$ . A method for determining a list  $A(\mathcal{P}, Z)$  of elements in  $\{0, 1, 2\}^{\mathcal{P}}$  adapted to sign determination for  $\mathcal{P}$  on  $Z$  from  $\text{SIGN}(\mathcal{P}, Z)$  has been given in Algorithm 10.12 (Family adapted to Sign Determination).

We are ready for describing the algorithm computing the Euler-Poincaré characteristic. We start with an algorithm for the Euler-Poincaré-query.

*Algorithm 13.4. [Euler-Poincaré-query]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$ , with  $Z = \text{Zer}(Q, \mathbb{R}^k)$ , a polynomial  $P \in D[X_1, \dots, X_k]$ .
- **Output:** the Euler-Poincaré-query

$$\text{EuQ}(P, Z) = \chi(P > 0, Z) - \chi(P < 0, Z).$$

- **Complexity:**  $d^{O(k)}$ , where  $d$  is a bound on the degree of  $Q$  and the degree of  $P$ .
- **Procedure:**
  - Introduce a new variable  $X_{k+1}$ , and let

$$\begin{aligned} Q_+ &= Q^2 + (P - X_{k+1}^2)^2, \\ Q_- &= Q^2 + (P + X_{k+1}^2)^2. \end{aligned}$$

Using Algorithm 12.25 compute  $\chi(\text{Zer}(Q_+, \mathbb{R}^{k+1}))$  and  $\chi(\text{Zer}(Q_-, \mathbb{R}^{k+1}))$ . Output

$$(\chi(\text{Zer}(Q_+, \mathbb{R}^{k+1})) - \chi(\text{Zer}(Q_-, \mathbb{R}^{k+1}))) / 2.$$

**Proof of correctness:** The algebraic set  $\text{Zer}(Q_+, \mathbb{R}^{k+1})$  is semi-algebraically homeomorphic to the disjoint union of two copies of the semi-algebraic set defined by  $(P > 0) \wedge (Q = 0)$ , and the algebraic set defined by  $(P = 0) \wedge (Q = 0)$ . Hence, using Proposition 6.57, we have that

$$2 \chi(P > 0, Z) = \chi(\text{Zer}(Q_+, \mathbb{R}^{k+1})) - \chi(\text{Zer}((Q, P), \mathbb{R}^k)).$$

Similarly, we have that

$$2 \chi(P < 0, Z) = \chi(\text{Zer}(Q_-, \mathbb{R}^{k+1})) - \chi(\text{Zer}((Q, P), \mathbb{R}^k)). \quad \square$$

**Complexity Analysis:** The complexity of the algorithm is  $d^{O(k)}$  using the complexity analysis of Algorithm 12.25.

When  $D = \mathbb{Z}$  and the bitsizes of the coefficients of  $P$  are bounded by  $\tau$ , the bitsizes of the intermediate computations and the output are bounded by  $O(k^2 d^2(\tau + \log_2(kd)))$ .  $\square$

We are now ready to describe an algorithm for computing the Euler-Poincaré characteristic of the realizations of sign conditions.

*Algorithm 13.5. [Euler-Poincaré Characteristic of Sign Conditions]*

- **Structure:** an ordered domain  $D$  contained in a real close field  $R$ .
- **Input:** an algebraic set  $Z = \text{Zer}(Q, R^k) \subset R^k$  and a finite list  $\mathcal{P}$  of polynomials in  $D[X_1, \dots, X_k]$ .
- **Output:** the list  $\chi(\mathcal{P}, Z)$ .
- **Complexity:**  $s^{k'+1} O(d)^k + s^{k'}((k' \log_2(s) + k \log_2(d)) d)^{O(k)}$ , where  $k'$  is the dimension of  $Z$ ,  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degree of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Let  $\mathcal{P} = P_1, \dots, P_s$ ,  $\mathcal{P}_i = P_1, \dots, P_i$ . Compute  $\text{SIGN}(\mathcal{P}, Z)$  using Algorithm 13.3 (Sampling on an Algebraic Set).
  - Determine a list  $A(\mathcal{P}, Z)$  adapted to sign determination for  $\mathcal{P}$  on  $Z$  using Algorithm 10.12 (Family adapted to Sign Determination).
  - Define  $A = A(\mathcal{P}, Z)$ ,  $M = M(\mathcal{P}^A, \text{SIGN}(\mathcal{P}, Z))$ .
  - Compute  $\text{EuQ}(\mathcal{P}^A, Z)$  using repeatedly Algorithm 13.4 (Euler-Poincaré-query).
  - Using

$$M \cdot \chi(\mathcal{P}, Z) = \text{EuQ}(\mathcal{P}^A, Z),$$

and the fact that  $M$  is invertible, compute  $\chi(\mathcal{P}, Z)$ .

**Proof of correctness:** Immediate from Proposition 13.41. □

**Complexity analysis:** By Proposition 7.35,

$$\#(\text{SIGN}(\mathcal{P}, Z)) \leq \sum_{0 \leq j \leq k'} \binom{s}{j} 4^j d (2d - 1)^{k-1} = s^{k'} O(d)^k.$$

The number of calls to to Algorithm 13.4 (Euler-Poincaré-query) is equal to  $\#(\text{SIGN}(\mathcal{P}, Z))$ . The calls to Algorithm 13.4 (Euler-Poincaré-query) are done for polynomials which are products of at most

$$\log_2(\#(\text{SIGN}(\mathcal{P}, Z))) = k' \log_2(s) + k (\log_2(d) + O(1)).$$

products of polynomials of the form  $P$  or  $P^2$ ,  $P \in \mathcal{P}$  by Proposition 10.71, hence of degree  $(k' \log_2(s) + k (\log_2(d) + O(1))) d$ . Using the complexity analysis of Algorithm 13.3 (Sampling on an Algebraic Set) and the complexity analysis of Algorithm 13.4 (Euler-Poincaré-query), the number of arithmetic operations is

$$s^{k'+1} O(d)^k + s^{k'} ((k' \log_2(s) + k \log_2(d)) d)^{O(k)}.$$

The algorithm also involves the inversion matrices of size  $s^{k'} O(d)^k$  with integer coefficients.

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ .  $\square$

### 13.5 Bibliographical Notes

Grigor'ev and Vorobjov [77] gave the first algorithm to solve the decision problem for the existential theory of the reals whose time complexity is singly exponential in the number of variables. Canny [37], Heintz, Roy, and Solerno [85], and Renegar [133] improved their result in several directions. Renegar's [133] algorithms solved the existential theory of the reals in time  $(s d)^{O(k)}$  (where  $d$  is the degree,  $k$  the number of variables, and  $s$  the number of polynomials). The first single exponential complexity computation for the Euler-Poincaré characteristic appears in [11].

The results presented in the three first sections are based on [13, 15]. The construction of the family  $\mathcal{L}_{k, k-k'}$  described in Section 13.3, is on the work of Chistov, Fournier, Gurvits, and Koiran [42]. In terms of algebraic complexity (the degree of the equations), they are similar to [133]. They are more precise in terms of combinatorial complexity (the dependence on the number of equations), particularly for the computation of the realizable sign conditions on a lower dimensional algebraic set.

## Quantifier Elimination

---

The principal problem we consider in this chapter is the quantifier elimination problem. This problem was already studied in Chapter 11, where we obtained doubly exponential complexity in the number of variables. On the other hand, we have seen in Chapter 13 an algorithm for the existential theory of the reals (which is to decide the truth or the falsity of a sentence with a single block of existential quantifiers) with complexity singly exponential in the number of variables (see Theorem 13.13). In this chapter, we pay special attention to the structure of the blocks of variables in a formula in order to take into account this block structure in the complexity estimates and improve the results obtained in Chapter 11.

If  $Z = (Z_1, \dots, Z_\ell)$ ,  $\Phi$  is a formula, and  $\text{Qu} \in \{\forall, \exists\}$ , we denote the formula  $(\text{Qu } Z_1) \dots (\text{Qu } Z_\ell) \Phi$  by the abbreviation  $(\text{Qu } Z) \Phi$ .

Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k, Y_1, \dots, Y_\ell]$  be finite, and let  $\Pi$  denote a partition of the list of variables  $X = (X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  is of size  $k_i$ ,  $1 \leq i \leq \omega$ ,  $\sum_{1 \leq i \leq \omega} k_i = k$ .

A  $(\mathcal{P}, \Pi)$ -formula  $\Phi(Y)$  is a formula of the form

$$\Phi(Y) = (\text{Qu}_1 X_{[1]}) \dots (\text{Qu}_\omega X_{[\omega]}) F(X, Y),$$

where  $\text{Qu}_i \in \{\forall, \exists\}$ ,  $Y = (Y_1, \dots, Y_\ell)$ , and  $F(X, Y)$  is a quantifier free  $\mathcal{P}$ -formula.

In Section 14.1, we describe an algorithm for solving the general decision problem, that is a procedure to decide the truth or falsity of a  $(\mathcal{P}, \Pi)$ -sentence. The key notion here is the tree of realizable sign conditions of a family of polynomials with respect to a block structure  $\Pi$  on the set of variables. This is a generalization of the set of realizable sign conditions, seen in Chapter 7, which corresponds to one single block of variables. It is also a generalization of the tree of cylindrical realizable sign conditions, seen in Chapter 11, which correspond to  $k$  blocks of one variable each. The basic idea of this algorithm is to perform parametrically the algorithm in Chapter 13, using the critical point method.

Section 14.2 is devoted to the more general problem of quantifier elimination for a  $(\mathcal{P}, \Pi)$ -formula.



Section 14.3 is devoted to a variant of Quantifier Elimination exploiting better the logical structure of the formula.

Finally, the block elimination technique is used to perform global optimization and compute the dimension of a semi-algebraic set in Section 14.4 and Section 14.5.

## 14.1 Algorithm for the General Decision Problem

We first study the general decision problem, which is to decide the truth or falsity of a  $(\mathcal{P}, \Pi)$ -sentence (which is a  $(\mathcal{P}, \Pi)$ -formula without free variables). In order to decide the truth or falsity of a sentence, we construct a certain tree of sign conditions adapted to the block structure  $\Pi$  of the sentence, which we define below.

The following definition generalizes the definition of the tree of cylindrical realizable sign conditions (Notation 11.7). The importance of this notion is that the truth or falsity of any  $(\mathcal{P}, \Pi)$ -sentence can be decided from  $\text{SIGN}_{\Pi}(\mathcal{P})$ .

**Notation 14.1. [Block realizable sign conditions]** Let  $\mathcal{P}$  be a set of  $s$  polynomials in  $k$  variables  $X_1, \dots, X_k$ , and let  $\Pi$  denote a partition of the list of variables  $X_1, \dots, X_k$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  is of size  $k_i$ , for  $1 \leq i \leq \omega$ ,  $\sum_{1 \leq i \leq \omega} k_i = k$ . Let  $\mathbb{R}^{[i]} = \mathbb{R}^{k_1 + \dots + k_i}$ , and let  $\pi_{[i]}$  be the projection from  $\mathbb{R}^{[i+1]}$  to  $\mathbb{R}^{[i]}$  forgetting the last  $k_{i+1}$ -coordinates. Note that  $\mathbb{R}^{[\omega]} = \mathbb{R}^k$ . By convention,  $\mathbb{R}^{[0]} = \{0\}$ .

We are going to define inductively the tree of realizable sign conditions of  $\mathcal{P}$  with respect to  $\Pi$ .

For  $z \in \mathbb{R}^{[\omega]}$ , let  $\text{SIGN}_{\Pi, \omega}(\mathcal{P})(z) = \text{sign}(\mathcal{P})(z)$ , where  $\text{sign}(\mathcal{P})(z)$  is the sign condition on  $\mathcal{P}$  mapping  $P \in \mathcal{P}$  to  $\text{sign}(P)(z) \in \{0, 1, -1\}$  (Notation 11.7).

For all  $i$ ,  $0 \leq i < \omega$ , and  $y \in \mathbb{R}^{[i]}$ , we inductively define,

$$\text{SIGN}_{\Pi, i}(\mathcal{P})(y) = \{\text{SIGN}_{\Pi, i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{[i+1]}, \pi_{[i]}(z) = y\}.$$

Finally, we define

$$\text{SIGN}_{\Pi}(\mathcal{P}) = \text{SIGN}_{\Pi, 0}(\mathcal{P})(0).$$

Note that  $\text{SIGN}_{\Pi}(\mathcal{P})$  is naturally equipped with a tree structure. We call  $\text{SIGN}_{\Pi}(\mathcal{P})$  the **tree of realizable sign conditions of  $\mathcal{P}$  with respect to  $\Pi$** .  $\square$

When there is only one block of variables, we recover  $\text{SIGN}(\mathcal{P})$  (Notation 7.29). When  $\Pi = \{X_1\}, \dots, \{X_k\}$ , we recover  $\text{CSIGN}(\mathcal{P})$  (Notation 11.7).

We will see that the truth or falsity of a  $(\mathcal{P}, \Pi)$ -sentence can be decided from the set  $\text{SIGN}_{\Pi}(\mathcal{P})$ . We first consider an example.

*Example 14.2.* Let  $P = X_1^2 + X_2^2 + X_3^2 - 1$ ,  $\mathcal{P} = \{P\}$ . Let  $\Pi$  consist of two blocks of variables, defined by  $X_{[1]} = X_1$  and  $X_{[2]} = \{X_2, X_3\}$ . Note that  $\pi_{[1]}$  projects  $\mathbb{R}^{[2]} = \mathbb{R}^3$  to  $\mathbb{R}^{[1]} = \mathbb{R}$  by forgetting the last two coordinates. We now determine  $\text{SIGN}_{\Pi}(\mathcal{P})$ .

For  $x \in \mathbb{R} = \mathbb{R}^{[1]}$ ,

$$\text{SIGN}_{\Pi,1}(\mathcal{P})(x) = \{\text{sign}(P(z)) \mid z \in \mathbb{R}^{[2]}, \pi_{[1]}(z) = x\}.$$

Thus

$$\text{SIGN}_{\Pi,1}(\mathcal{P})(x) = \begin{cases} \{0, 1, -1\} & \text{if } x \in (-1, 1) \\ \{0, 1\} & \text{if } x \in \{-1, 1\} \\ \{1\} & \text{otherwise.} \end{cases}$$

Finally,

$$\text{SIGN}_{\Pi}(\mathcal{P}) = \{\text{SIGN}_{\Pi,1}(\mathcal{P})(x) \mid x \in \mathbb{R}\}.$$

Thus

$$\text{SIGN}_{\Pi}(\mathcal{P}) = \{\{1\}, \{0, 1\}, \{0, 1, -1\}\}.$$

This means that there are three cases:

- there are values of  $x_1$  for which the only sign taken by  $P(x_1, x_2, x_3)$  when  $(x_2, x_3)$  varies in  $\mathbb{R}^2$  is 1,
- there are values of  $x_1$  for which the only sign taken by  $P(x_1, x_2, x_3)$  when  $(x_2, x_3)$  varies in  $\mathbb{R}^2$  are 0 and 1,
- there are values of  $x_1$  for which the signs taken by  $P(x_1, x_2, x_3)$  when  $(x_2, x_3)$  varies in  $\mathbb{R}^2$  are 0, 1 and  $-1$ ,
- and these exhaust all choices of  $x_1 \in \mathbb{R}$ .

So, the sentence  $(\forall X_1) (\exists (X_2, X_3)) X_1^2 + X_2^2 + X_3^2 - 1 > 0$  is certainly true.

Since there are values of  $x_1$  for which the only sign taken by  $P(x_1, x_2, x_3)$  for every  $(x_2, x_3) \in \mathbb{R}^2$  is 1 it is equally clear that the sentence  $(\exists X_1) (\forall (X_2, X_3)) X_1^2 + X_2^2 + X_3^2 - 1 > 0$  is true.

On the other hand, the sentence  $(\forall X_1) (\exists (X_2, X_3)) X_1^2 + X_2^2 + X_3^2 - 1 = 0$  is false: there are values of  $x_1$  for which the only sign taken by  $P(x_1, x_2, x_3)$  is 1.

This differs from what was done in Example 11.10 in that here we do not decompose the  $(X_2, X_3)$  space: this is because the variables  $\{X_2, X_3\}$  belong to the same block of quantifiers. So the information provided by  $\text{SIGN}_{\Pi}(\mathcal{P})$  is weaker than the information provided by  $\text{CSIGN}(\mathcal{P})$  (Notation 11.7). Note that  $\text{SIGN}_{\Pi}(\mathcal{P})$  does not provide the information necessary to decide the truth or falsity of the sentence

$$\Phi = (\exists X_1) (\forall X_2) (\exists X_3) X_1^2 + X_2^2 + X_3^2 - 1 = 0$$

since we do not have information for the corresponding block structure, while we have able to decide that  $\Phi$  is false using

$\text{CSIGN}(\mathcal{P}) = \{\{\{\{1\}, \{0, 1\}, \{0, 1, -1\}\}, \{\{1\}, \{0, 1\}\}, \{\{1\}\}\}$  in Example 11.16.

If we take  $\mathcal{Q} = \{X_1 - X_3^2\}$ , it is easy to check that

$$\text{SIGN}_{\Pi}(\mathcal{Q}) = \{\{1\}, \{0, 1\}, \{0, 1, -1\}\} = \text{SIGN}_{\Pi}(\mathcal{P}).$$

On the other hand we can determine

$$\text{CSIGN}(\mathcal{Q}) = \{\{\{1\}\}, \{\{0, 1\}\}, \{\{0, 1, -1\}\}\}$$

and notice that

$$\text{CSIGN}(\mathcal{Q}) \neq \text{CSIGN}(\mathcal{P}).$$

Using  $\text{CSIGN}(\mathcal{Q})$ , we can check that the sentence

$$\Phi' = (\exists X_1) (\forall X_2) (\exists X_3) X_1 - X_3^2 = 0$$

while the corresponding  $\Phi$ , discussed above, is false.  $\square$

We use again Notation 11.12.

**Proposition 14.3.** *The  $(\mathcal{P}, \Pi)$ -sentence*

$$(\text{Qu}_1 X_{[1]}) (\text{Qu}_2 X_{[2]}) \dots (\text{Qu}_\omega X_{[\omega]}) F(X),$$

*is true if and only if*

$$\text{Qu}_1 \sigma_1 \in \text{SIGN}_{\Pi}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_\omega \sigma_\omega \in \sigma_{\omega-1} \quad F^*(\sigma_\omega).$$

**Proof:** The proof is by induction on the number  $\omega$  of blocks of quantifiers, starting from the one outside.

Since  $(\forall X) \Phi$  is equivalent to  $\neg(\exists X) \neg\Phi$ , we can suppose without loss of generality that  $\text{Qu}_1$  is  $\exists$ .

The claim is certainly true when there is one block of existential quantifiers, by definition of  $\text{sign}(\mathcal{P})$ .

Suppose that

$$(\exists X_{[1]}) (\text{Qu}_2 X_{[2]}) \dots (\text{Qu}_\omega X_{[\omega]}) F(X)$$

is true, and choose  $a_{[1]} \in \mathbb{R}^{k_1}$  such that

$$(\text{Qu}_2 X_{[2]}) \dots (\text{Qu}_\omega X_{[\omega]}) F(a_{[1]}, X_{[2]}, \dots, X_{[\omega]})$$

is true. Note that if  $\mathcal{P}_{a_{[1]}}$  is the set of polynomials obtained by substituting  $a_{[1]} \in \mathbb{R}^{k_1}$  for  $X_{[1]}$  in  $\mathcal{P}$  and  $\Pi' = X_{[2]}, \dots, X_{[\omega]}$ ,

$$\text{SIGN}_{\Pi,1}(\mathcal{P})(a_{[1]}) = \text{SIGN}_{\Pi'}(\mathcal{P}_{a_{[1]}}).$$

By induction hypothesis,

$$\text{Qu}_2 \sigma_2 \in \text{SIGN}_{\Pi'}(\mathcal{P}_{a_{[1]}}) \dots \text{Qu}_\omega \sigma_\omega \in \sigma_{\omega-1} \quad F^*(\sigma_\omega)$$

is true. So taking  $\sigma_1 = \text{SIGN}_{\Pi,1}(\mathcal{P})(a_{[1]}) = \text{SIGN}_{\Pi'}(\mathcal{P}_{a_{[1]}}) \in \text{SIGN}_{\Pi}(\mathcal{P})$ ,

$$\exists \sigma_1 \in \text{SIGN}_{\Pi}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_{\omega} \sigma_{\omega} \in \sigma_{\omega-1} \quad F^*(\sigma_{\omega})$$

is true.

Conversely, suppose

$$\exists \sigma_1 \in \text{SIGN}_{\Pi}(\mathcal{P}) \quad \text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_{\omega} \sigma_{\omega} \in \sigma_{\omega-1} \quad F^*(\sigma_{\omega})$$

is true and choose  $\sigma_1 \in \text{SIGN}_{\Pi}(\mathcal{P})$  such that

$$\text{Qu}_2 \sigma_2 \in \sigma_1 \dots \text{Qu}_{\omega} \sigma_{\omega} \in \sigma_{\omega-1} \quad F^*(\sigma_{\omega})$$

is true. By definition of  $\text{SIGN}_{\Pi}(\mathcal{P})$ ,  $\sigma_1 = \text{SIGN}_{\Pi'}(\mathcal{P})(a_{[1]})$  for some  $a_{[1]} \in \mathbb{R}^{k_1}$ , and hence

$$\text{Qu}_2 \sigma_2 \in \text{SIGN}_{\Pi'}(\mathcal{P}_{a_{[1]}}) \dots \text{Qu}_{\omega} \sigma_{\omega} \in \sigma_{\omega-1} \quad F^*(\sigma_{\omega})$$

is true. By induction hypothesis,

$$(\text{Qu}_2 X_{[2]}) \dots (\text{Qu}_{\omega} X_{[\omega]}) F(a_{[1]}, X_{[2]}, \dots, X_{[\omega]})$$

is true. Thus

$$(\exists X_{[1]}) (\text{Qu}_2 X_{[2]}) \dots (\text{Qu}_{\omega} X_{[\omega]}) F(X)$$

is true. □

In the cylindrical situation studied in Chapter 11,  $\text{CSIGN}(\mathcal{P})$  was obtained from a cylindrical set of sample points of a cylindrical decomposition adapted to  $\mathcal{P}$ . We generalize this approach to a general block structure.

A  $\Pi$ -set  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_{\omega}$  is such that  $\mathcal{A}_i$  is a finite set contained in  $\mathbb{R}^{[i]}$  and  $\pi_{[i]}(\mathcal{A}_{i+1}) = \mathcal{A}_i$ .

We define inductively the **tree of realizable sign conditions of  $\mathcal{P}$  for  $\mathcal{A}$  with respect to  $\Pi$** ,  $\text{SIGN}_{\Pi}(\mathcal{P}, \mathcal{A})$ , as follows:

- For  $z \in \mathcal{A}_{\omega}$ , let  $\text{SIGN}_{\Pi,\omega}(\mathcal{P})(z) = \text{sign}(\mathcal{P})(z)$ , where  $\text{sign}(\mathcal{P})(z)$  is the sign condition on  $\mathcal{P}$  mapping  $P \in \mathcal{P}$  to  $\text{sign}(P)(z) \in \{0, 1, -1\}$  (Notation 11.7).
- For all  $i$ ,  $1 \leq i < \omega$ , and all  $y \in \mathcal{A}_i$ , we inductively define,

$$\text{SIGN}_{\Pi,i}(\mathcal{P}, \mathcal{A})(y) = \{\text{SIGN}_{\Pi,i+1}(\mathcal{P}, \mathcal{A})(z) \mid z \in \mathcal{A}_{i+1}, \pi_{[i]}(z) = y\}.$$

Finally, we define

$$\text{SIGN}_{\Pi}(\mathcal{P}, \mathcal{A}) = \text{SIGN}_{\Pi,0}(\mathcal{P}, \mathcal{A})(0).$$

Note that  $\text{SIGN}_{\Pi}(\mathcal{P}) = \text{SIGN}_{\Pi}(\mathcal{P}, \mathbb{R}^k)$ . Note also that  $\text{SIGN}_{\Pi}(\mathcal{P}, \mathcal{A})$  is a subtree of  $\text{SIGN}_{\Pi}(\mathcal{P})$ .

A  $\Pi$ -set of sample points for  $\mathcal{P}$  is a  $\Pi$ -set  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_{\omega}$  such that

$$\text{SIGN}_{\Pi}(\mathcal{P}, \mathcal{A}) = \text{SIGN}_{\Pi}(\mathcal{P}).$$

A  $\Pi$ -partition adapted to  $\mathcal{P}$  is given by  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_\omega$ , where  $\mathcal{S}_i$  is a partition of  $\mathbb{R}^{[i]}$  into a finite number of semi-algebraically connected semi-algebraic sets such that for every  $S \in \mathcal{S}_{i+1}$ ,  $\pi_{[i]}(S) \in \mathcal{S}_i$ , and such that every  $S \in \mathcal{S}_\omega$  is  $\mathcal{P}$ -invariant. A  $\Pi$ -set of sample points for a  $\Pi$ -partition  $\mathcal{S}$  is a  $\Pi$ -set  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_\omega$  such that

- for every  $i$ ,  $1 \leq i \leq \omega$ ,  $\mathcal{A}_i$  intersects every  $S \in \mathcal{S}_i$ ,
- for every  $i$ ,  $1 \leq i \leq \omega - 1$ ,  $\pi_{[i]}(\mathcal{A}_{i+1}) = \mathcal{A}_i$ .

Note that the partition of  $\mathbb{R}^k$  by the semi-algebraically connected components of realizable sign conditions of  $\mathcal{P}$  is a  $\Pi$ -partition with the block structure  $\Pi = \{X_1, \dots, X_k\}$  (i.e. with a single block), and a set of sample points for  $\mathcal{P}$  is a  $\Pi$ -set of sample points for this block structure. Note also that a cylindrical decomposition  $\mathcal{S}$  adapted to  $\mathcal{P}$  is a  $\Pi$ -partition for the block structure  $X_1, \dots, X_k$  ( $k$ -blocks of one variable) and a cylindrical set of sample points for  $\mathcal{S}$  is a  $\Pi$ -set of sample points for  $\mathcal{S}$  for this block structure.

We are going to prove a result generalizing Proposition 11.9 to the case of a general block structure.

**Proposition 14.4.** *Let  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_\omega$  be a  $\Pi$ -partition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$  and  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_\omega$  be a  $\Pi$ -set of sample points for  $\mathcal{S}$ . Then  $\mathcal{A}$  is a  $\Pi$ -set of sample points for  $\mathcal{P}$ .*

The proof is similar to the proof of Proposition 11.9 and uses the following generalization of Proposition 11.11.

**Proposition 14.5.** *Let  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_\omega$  be a  $\Pi$ -partition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ . For every  $1 \leq i \leq \omega$  and every  $S \in \mathcal{S}_i$ ,  $\text{SIGN}_{\Pi,i}(y)$  is constant as  $y$  varies in  $S$ .*

**Proof:** The proof is by induction on  $\omega - i$ .

If  $i = \omega$ , the claim is true since the sign of every  $P \in \mathcal{P}$  is fixed on  $S \in \mathcal{S}_\omega$ .

Suppose that the claim is true for  $i + 1$  and consider  $S \in \mathcal{S}_i$ . Let  $T_1, \dots, T_\ell$  be the elements of  $\mathcal{S}_{i+1}$  such that  $\pi_{[i]}(T_j) = S$ . By induction hypothesis,  $\text{SIGN}_{\Pi,i+1}(\mathcal{P})(z)$  is constant as  $z$  varies in  $T_j$ . Since  $\mathcal{S}$  is a  $\Pi$ -partition,  $\bigcup_{j=1}^\ell T_j = S \times \mathbb{R}^{k_{i+1}}$ . Thus

$$\text{SIGN}_{\Pi,i}(\mathcal{P})(y) = \{\text{SIGN}_{\Pi,i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{[i+1]}, \pi_{[i]}(z) = y\}$$

is constant as  $y$  varies in  $S$ . □

**Proof of Proposition 14.4:** Let  $\mathcal{A}_0 = \{0\}$ . We are going to prove that for every  $y \in \mathcal{A}_i$ ,

$$\text{SIGN}_{\Pi,i}(\mathcal{P})(y) = \text{SIGN}_{\Pi,i}(\mathcal{P}, \mathcal{A})(y).$$

The proof is by induction on  $\omega - i$ .

If  $i = \omega$ , the claim is true since  $\mathcal{A}_\omega$  meets every element of  $\mathcal{S}_\omega$ .

Suppose that the claim is true for  $i + 1$  and consider  $y \in \mathcal{A}_i$ . Let  $S$  be the element of  $\mathcal{S}_i$  containing  $y$ , and let  $T_1, \dots, T_\ell$  be the elements of  $\mathcal{S}_{i+1}$  such that  $\pi_{[i]}(T_j) = S$ . Denote by  $z_j$  a point of  $T_j \cap \mathcal{A}_{i+1}$  such that  $\pi_{[i]}(z_j) = y$ . By induction hypothesis,

$$\text{SIGN}_{\Pi, i+1}(\mathcal{P})(z_j) = \text{SIGN}_{\Pi, i+1}(\mathcal{P}, \mathcal{A})(z_j).$$

Since  $T_1 \cup \dots \cup T_\ell = S \times \mathbb{R}^{k_{i+1}}$  and  $\text{SIGN}_{\Pi, i+1}(\mathcal{P})(z)$  does not change as  $z$  varies over  $T_j$ ,

$$\begin{aligned} \text{SIGN}_{\Pi, i}(\mathcal{P})(y) &= \{\text{SIGN}_{\Pi, i+1}(\mathcal{P})(z) \mid z \in \mathbb{R}^{[i+1]}, \pi_{[i]}(z) = y\} \\ &= \{\text{SIGN}_{\Pi, i+1}(\mathcal{P}, \mathcal{A})(z) \mid z \in \mathcal{A}_{i+1}, \pi_{[i]}(z) = y\} \\ &= \text{SIGN}_{\pi, i}(\mathcal{P}, \mathcal{A})(y). \end{aligned}$$

□

We now construct a  $\Pi$ -partition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ , generalizing Theorem 5.6. Note that a cylindrical decomposition adapted to  $\mathcal{P}$  gives a  $\Pi$ -partition of  $\mathbb{R}^k$  adapted to  $\mathcal{P}$ , so the issue here is not an existence theorem similar to Theorem 5.6 but rather a complexity result taking into account the block structure  $\Pi$ . The construction of a cylindrical decomposition adapted to  $\mathcal{P}$  in Chapter 5 and Chapter 11 was based on a recursive call to an Elimination procedure eliminating one variable (see Algorithm 11.1 (Subresultant Elimination)). In the general block structure context, we define a Block Elimination procedure which replaces a block of variables by one single variable and computes parametrized univariate representations, giving in a parametric way sample points for every non-empty sign condition. Finally we eliminate this variable.

*Algorithm 14.1.* **[Block Elimination]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a block of variables  $X = (X_1, \dots, X_k)$  and a set of polynomials

$$\mathcal{P}(Y) \subset D[Y_1, \dots, Y_\ell, X_1, \dots, X_k].$$

- **Output:**
  - a set  $\text{BElim}_X(\mathcal{P}) \subset D[Y]$  such that  $\text{SIGN}(\mathcal{P}(y, X_1, \dots, X_k))$  (Notation 7.29) is fixed as  $y$  varies over a semi-algebraically connected component of a realizable sign condition of  $\text{BElim}_X(\mathcal{P})$ ,
  - a set  $\text{UR}_X(\mathcal{P})$  of parametrized univariate representations of the form

$$u(Y, \varepsilon, \delta) = (f, g_0, \dots, g_k),$$

where  $f, g_i \in D[Y, \varepsilon, \delta][T]$ . For any point  $y \in \mathbb{R}^\ell$ , denoting by  $\text{UR}_X(\mathcal{P})(y)$  the subset of  $\text{UR}_X(\mathcal{P})$  such that  $f(y, T)$  and  $g_0(y, T)$  are coprime, the points associated to the univariate representations in  $\text{UR}_X(\mathcal{P})(y)$  intersect every semi-algebraically connected component of every realizable sign condition of the set  $\mathcal{P}(y)$  in  $\mathbb{R}(\varepsilon, \delta)^k$ .

- **Complexity:**  $s^{k+1}d^{O(\ell k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on their degree.
- **Procedure:**
  - Initialize  $\text{UR}_X(\mathcal{P})$  to the empty set.
  - Take as  $d'$  the smallest even natural number  $> d$ .
  - Define

$$P_i^* = \{(1 - \delta)P_i + \delta H_k(d', i), (1 - \delta)P_i - \delta H_k(d', i), \\ (1 - \delta)P_i + \delta \gamma H_k(d', i), (1 - \delta)P_i - \delta \gamma H_k(d', i)\}$$

$$\mathcal{P}^* = \{P_1^*, \dots, P_s^*\}$$

for  $0 \leq i \leq s$  using Notation 13.4.

- For every subset  $\mathcal{Q}$  of  $j \leq k$  polynomials  $Q_{i_1} \in P_{i_1}^*, \dots, Q_{i_j} \in P_{i_j}^*$ ,
  - let  $Q = Q_{i_1}^2 + \dots + Q_{i_j}^2 + (\varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2$ .
  - Take for  $i = 1, \dots, k$ ,  $\bar{d}_i$  the smallest even natural number  $> \deg(Q)$ ,  $\bar{d}_{k+1} = 6$ ,  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k, \bar{d}_{k+1})$ , and  $c = \varepsilon$ .
  - Perform Algorithm 12.10 (Parametrized Multiplication Table) with input  $\overline{\text{Cr}}(Q, \zeta)$  (using Notation 12.46). Output  $\mathcal{M}$ .
  - Perform Algorithm 12.15 (Parametrized Limit of Bounded Points) with input  $\gamma, \zeta, \overline{\text{Cr}}(Q, \zeta)$ , and  $\mathcal{M}$ . Add the parametrized univariate representations (belonging to  $\text{D}[Y, \varepsilon, \delta][T]^{k+2}$ ) output to  $\text{UR}_X(\mathcal{P})$ .
- For every  $v = (f, g_0, \dots, g_k) \in \text{UR}_X(\mathcal{P})$ , consider the family of univariate polynomials  $\mathcal{F}_v$  consisting of  $f$ , its derivatives with respect to  $T$ , and  $P_v$  (see Notation 13.8), for every  $P \in \mathcal{P}$ . Compute  $\text{RElim}_T(f, \mathcal{F}_v)$  using Algorithm 11.19 (Restricted Elimination). Denote by  $\mathcal{B}_v$  the family of polynomials in  $Y$  that are the coefficients of the polynomials in

$$\text{RElim}_T(\mathcal{F}_v) \subset \text{D}[Y, \varepsilon, \delta].$$

- Define  $\text{BElim}_X(\mathcal{P})$  to be the union of the sets  $\mathcal{B}_v \subset \text{D}[Y]$  for every  $v \in \text{UR}_X(\mathcal{P})$ .
- Output  $\text{BElim}_X(\mathcal{P})$  and  $\text{UR}_X(\mathcal{P})$ .

The proof of correctness of Algorithm 14.1 (Block Elimination) uses the following results, which describe how to get rid of infinitesimal quantities.

**Notation 14.6.** Let  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_m$  be variables and consider the real closed field  $\mathbb{R}\langle \varepsilon_1, \varepsilon_2, \dots, \varepsilon_m \rangle$ . Let  $S \subset \mathbb{R}\langle \varepsilon_1, \dots, \varepsilon_m \rangle^k$  be a semi-algebraic set defined by a quantifier-free  $\mathcal{P}$ -formula  $\Phi$  with  $\mathcal{P} \subset \text{D}[\varepsilon_1, \dots, \varepsilon_m, X_1, \dots, X_k]$ . Let  $P \in \mathcal{P}$ . We write  $P$  as a polynomial in  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$  and order the monomials with the order induced by the order  $<_\varepsilon$  on  $\mathbb{R}\langle \varepsilon_1, \varepsilon_2, \dots, \varepsilon_m \rangle$  with  $\varepsilon_1 >_\varepsilon \dots >_\varepsilon \varepsilon_m$ . Let

$$P = P_0\varepsilon^{\alpha_0} + P_1\varepsilon^{\alpha_1} + \dots + P_m\varepsilon^{\alpha_m},$$

where,  $P_i \in \text{D}[X_1, \dots, X_k]$ ,  $\alpha_i \in \mathbb{N}^\ell$ , and  $\varepsilon^{\alpha_i} >_\varepsilon \varepsilon^{\alpha_{i+1}}$ , for  $0 \leq i \leq m$ .

Define

$$\text{Remo}_\varepsilon(P=0) = \bigwedge_{i=0}^m (P_i=0),$$

$$\text{Remo}_\varepsilon(P>0) = (P_0>0) \vee (P_0=0 \wedge P_1>0) \vee \dots \vee \left( \bigwedge_{i=0}^{m-1} P_i=0 \wedge P_m>0 \right),$$

$$\text{Remo}_\varepsilon(P<0) = (P_0<0) \vee (P_0=0 \wedge P_1<0) \vee \dots \vee \left( \bigwedge_{i=0}^{m-1} P_i=0 \wedge P_m<0 \right).$$

Let  $\text{Remo}_\varepsilon(\Phi)$  be the formula obtained from  $\Phi$  by replacing every atom,  $P=0$ ,  $P>0$ , or  $P<0$  in  $\Phi$  by the corresponding formula

$$\text{Remo}_\varepsilon(P=0), \text{Remo}_\varepsilon(P>0), \text{Remo}_\varepsilon(P<0). \quad \square$$

**Proposition 14.7.** *Let  $S \subset \mathbb{R}\langle \epsilon_1, \dots, \epsilon_m \rangle^k$  be a semi-algebraic set defined by a quantifier-free  $\mathcal{P}$ -formula  $\Phi$  with  $\mathcal{P} \subset \mathbb{D}[\epsilon_1, \dots, \epsilon_m, X_1, \dots, X_k]$ . Let  $S' \subset \mathbb{R}^k$  be the semi-algebraic set defined by  $\text{Remo}_\varepsilon(\Phi)$ . Then,  $S' = S \cap \mathbb{R}^k$ .*

**Proof:** Let  $x \in \mathbb{R}^k$  satisfy  $\text{Remo}_\varepsilon(\Phi)$ . It is clear by construction that  $x$  also satisfies  $\Phi$ . Conversely, if  $x \in S \cap \mathbb{R}^k$ , then for any polynomial

$$P \in \mathbb{D}[\epsilon_1, \dots, \epsilon_m, X_1, \dots, X_k],$$

the sign of  $P(x)$  is determined by the sign of the coefficient of the biggest monomial in the lexicographical ordering, when  $P(x)$  is expressed as a polynomial in  $\epsilon_1, \dots, \epsilon_m$ . This immediately implies that  $x$  satisfies the formula  $\text{Remo}_\varepsilon(\Phi)$ .  $\square$

**Proof of correctness of Algorithm 14.1:** The result follows from the correctness of Algorithm 13.1 (Computing Realizable Sign Conditions) and Algorithm 11.19 (Restricted Elimination). Consider a semi-algebraically connected component  $S$  of a realizable sign condition on  $\text{BELim}_X(\mathcal{P})$ . Then, the following remain invariant as  $y$  varies over  $S$ : the set  $\text{UR}_X(\mathcal{P})(y)$ , for every

$$u(Y, \varepsilon, \delta) = (f(Y, \varepsilon, \delta, T), g_0(Y, \varepsilon, \delta, T), \dots, g_k(Y, \varepsilon, \delta, T)) \in \text{UR}_X(\mathcal{P})(y),$$

the number of roots of  $f(y, \varepsilon, \delta, T)$  in  $\mathbb{R}\langle \varepsilon, \delta \rangle$  and their Thom encodings, as well as the number of roots in  $\mathbb{R}\langle \varepsilon, \delta \rangle$  that are common to  $f(y, \varepsilon, \delta, T)$  and  $P_u(y, \varepsilon, \delta, T)$  for all  $P \in \mathcal{P}$ . These are consequences of the properties of  $\text{RElim}_T$  (see Algorithm 11.19 (Restricted Elimination)). It is finally clear that  $\text{SIGN}(\mathcal{P}(y, X))$  is constant as  $y$  varies in a semi-algebraically connected component  $S$  of a realizable sign condition on  $\text{BELim}_X(\mathcal{P})$ , using Proposition 14.7.  $\square$

**Complexity analysis of Algorithm 14.1:** The number of arithmetic operations in  $\mathbb{D}[Y, \varepsilon, \delta, \gamma, \zeta]$  for computing

$$\text{UR}_X(\mathcal{P}) \subset \mathbb{D}[Y, \varepsilon, \delta][T]$$



is  $\sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)}$ , using the complexity analysis of Algorithm 13.1 (Computing Realizable Sign Conditions). The degrees of the polynomials in  $T$  generated in this process are bounded by  $O(d)^k$  (independent of  $\ell$ ), and the degree in the variables  $Y$  as well as in the variables  $\varepsilon$  and  $\delta$  is  $d^{O(k)}$ , using the complexity analysis of Algorithm 12.10 (Parametrized Multiplication Table) and Algorithm 12.15 (Parametrized Limit of Bounded Points).

The complexity in  $D$  for computing  $\text{UR}_X(\mathcal{P})$  is  $\sum_{j \leq k} 4^j \binom{s}{j} d^{O(\ell k)}$ , using the complexity analysis of Algorithm 8.4 (Addition of multivariate polynomials) and Algorithm 8.5 (Multiplication of multivariate polynomials).

Using the complexity of Algorithm 11.19 (Restricted Elimination), the size of the set  $\text{BElim}_X(\mathcal{P})$  is  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(k)} = s^{k+1} d^{O(k)}$ , and the degrees of the elements of  $\text{BElim}_X(\mathcal{P})$  is  $d^{O(k)}$ .

The complexity in  $D$  is finally  $s \sum_{j \leq k} 4^j \binom{s}{j} d^{O(\ell k)} = s^{k+1} d^{O(\ell k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ .  $\square$

We construct a  $\Pi$ -partition adapted to  $\mathcal{P}$  using recursive calls to Algorithm 14.1 (Block Elimination).

**Notation 14.8.** Defining  $\text{B}_{\Pi, \omega}(\mathcal{P}) = \mathcal{P}$ , we denote, for  $1 \leq i \leq \omega - 1$ ,

$$\text{B}_{\Pi, i}(\mathcal{P}) = \text{BElim}_{X_{[i+1]}}(\text{B}_{\Pi, i+1}(\mathcal{P})),$$

so that  $\text{B}_{\Pi, i}(\mathcal{P}) \subset \mathbb{R}[X_{[1]}, \dots, X_{[i]}]$ .  $\square$

For every  $i$ ,  $1 \leq i \leq \omega$ , let  $\mathcal{S}_i$  be the set of semi-algebraically connected components of non-empty realizations of sign conditions on  $\bigcup_{j=1}^i \text{B}_{\Pi, i}(\mathcal{P})$ .

The following proposition follows clearly from the correctness of Algorithm 14.1 (Block Elimination).

**Proposition 14.9.** *The list  $\mathcal{S} = \mathcal{S}_1, \dots, \mathcal{S}_\omega$  is a  $\Pi$ -partition adapted to  $\mathcal{P}$ .*

In order to describe a  $\Pi$ -set of sample points for  $\mathcal{S}$ , we are going to use the parametrized univariate representations computed in Algorithm 14.1 (Block Elimination).

**Notation 14.10.** *Note that for every  $i = \omega - 1, \dots, 0$ ,*

$$\text{UR}_{\Pi, i}(\mathcal{P}) = \text{UR}_{X_{[i+1]}} \text{B}_{\Pi, i+1}(\mathcal{P}).$$

The elements of  $\text{UR}_{\Pi, i}(\mathcal{P})$  are parametrized univariate representations in the variable  $T_{i+1}$ , contained in  $D[X_{[1]}, \dots, X_{[i]}, \varepsilon_{i+1}, \delta_{i+1}][T_{i+1}]^{k_{i+1}+2}$ . Let

$$u = (u_0, \dots, u_{\omega-1}) \in \mathcal{U} = \prod_{i=0}^{\omega-1} \text{UR}_{\Pi, i}(\mathcal{P}),$$

with

$$u_{i-1} = (f^{[i]}, g_0^{[i]}, g_1^{[i]}, \dots, g_{k_i}^{[i]}).$$

For a polynomial  $P(X_{[1]}, \dots, X_{[i]})$ , let  $P_{u,i..j}(X_{[1]}, \dots, X_{[j-1]}, T_j, \dots, T_i)$  denote the polynomial obtained by successively replacing the blocks of variables  $X_{[\ell]}$ , with the rational fractions associated with the tuple  $u_{\ell-1}$  (using Notation ?), for  $\ell$  from  $i$  to  $j$ . Denoting  $P_{u,i}(T_1, \dots, T_i) = P_{u,i..1}(T_1, \dots, T_i)$ , define

$$\begin{aligned} \mathcal{T}_{u,i} &= (f^{[1]}(T_1), f_{u,1}^{[2]}(T_1, T_2), \dots, f_{u,i-1}^{[i]}(T_1, T_2, \dots, T_i)), \\ \mathcal{T}_u &= (f^{[1]}(T_1), f_{u,1}^{[2]}(T_1, T_2), \dots, f_{u,\omega-1}^{[\omega]}(T_1, T_2, \dots, T_\omega)), \\ \bar{u}_{i-1} &= (f_{u,i-1}^{[i]}, g_{0\ u,i-1}^{[i]}, g_{1\ u,i-1}^{[i]}, \dots, g_{k_{iu,i-1}}^{[i]}) \end{aligned}$$

Note that  $\bar{u}_{i-1}$  are univariate representations contained in

$$D[T_1, \dots, T_{i-1}, \varepsilon_1, \delta_1, \dots, \varepsilon_i, \delta_i][T_i]^{k_i+2}.$$

For  $u \in \mathcal{U}$  and  $t_\sigma \in \text{Zer}(\mathcal{T}_u, R\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle)$ , with Thom encoding  $\sigma$  let  $x_{u,\sigma,i} \in R\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle^{[i]}$  be the point obtained by substituting  $t_\sigma$  in the rational functions associated to  $\bar{u}_{j-1}$ ,  $j \leq i$ . Let  $\mathcal{A}_i$  be the set of points  $x_{u,\sigma,i} \in R\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle^{[i]}$  obtained by considering all  $u \in \mathcal{U}$  and  $t_\sigma \in \text{Zer}(\mathcal{T}_u, R\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle)$ . Then  $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_\omega$  is a  $\Pi$ -set, specified by  $\mathcal{V}$  where the elements of  $\mathcal{V}$  are pairs of an element  $u \in \mathcal{U}$  and a Thom encoding  $\sigma$  of an element of  $\text{Zer}(\mathcal{T}_u, R\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle)$ .  $\square$

The correctness of Algorithm 14.1 (Block Elimination) implies the following proposition.

**Proposition 14.11.** *The set  $\mathcal{A}$  is a  $\Pi$ -set of sample points for  $\mathcal{P}$ .*

Thus, in order to construct the set  $\text{SIGN}_\Pi(\mathcal{P})$ , it suffices to compute the signs of  $\mathcal{P}_{u,\omega}$  at the zeros of  $\mathcal{T}_u$ ,  $u \in \mathcal{U}$ .

The algorithm is as follows, using the notation of Algorithm 14.1 (Block Elimination):

*Algorithm 14.2.* **[Block Structured Signs]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:** a set  $\mathcal{P} \subset R[X_1, \dots, X_k]$ , and a partition,  $\Pi$ , of the variables  $X_1, \dots, X_k$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ .
- **Output:** the tree  $\text{SIGN}_\Pi(\mathcal{P})$  of realizable sign conditions of  $\mathcal{P}$  with respect to  $\Pi$ .
- **Complexity:**  $s^{(k_\omega+1)\dots(k_1+1)} d^{O(k_\omega)\dots O(k_1)}$ , where  $s$  is bound on the number of elements of  $\mathcal{P}$ ,  $d$  is a bound on their degree, and  $k_{[i]}$  is the number of elements of  $X_{[i]}$ .
- **Procedure:**
  - Initialize  $B_{\Pi,\omega}(\mathcal{P}) := \mathcal{P}$ .
  - Block Elimination Phase: Compute

$$B_{\Pi,i}(\mathcal{P}) = \text{BElim}_{X_{[i+1]}}(B_{\Pi,i+1}(\mathcal{P})),$$

for  $1 \leq i \leq \omega - 1$ , applying repeatedly  $\text{BElim}_{X_{[i+1]}}$ , using Algorithm 14.1 (Block Elimination). Define  $\text{B}_{\Pi,0}(\mathcal{P}) = \{1\}$ . Compute  $\text{UR}_{\Pi,i}(\mathcal{P})$ , for every  $i = \omega - 1, \dots, 0$ , using Algorithm 14.1 (Block Elimination). The elements of  $\text{UR}_{\Pi,i}(\mathcal{P})$  are parametrized univariate representations in the variable  $T_{i+1}$ , contained in

$$D[X_{[1]}, \dots, X_{[i]}, \varepsilon_{i+1}, \delta_{i+1}][T_{i+1}]^{k_{i+1}+2}.$$

- Substitution Phase: Compute the set of pairs  $\{(\mathcal{T}_u, \mathcal{P}_{u,\omega}) \mid u \in \mathcal{U}\}$ , using their definition in Notation 14.10.
- Sign Determination Phase: Compute the signs of the set of the polynomials in  $\mathcal{P}_{u,\omega}$  on  $\text{Zer}(\mathcal{T}_u, \mathbf{R}\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle^\omega)$  using Algorithm 12.19 (Zero-dimensional Sign Determination).
- Construct the set  $\text{SIGN}_{\Pi}(\mathcal{P})$  from these signs.

**Proof of correctness:** The correctness of the algorithm follows from Proposition 14.11.  $\square$

**Complexity analysis:** Using the complexity of Algorithm 14.1 (Block Elimination), the degrees and number of the parametrized univariate representations in  $\text{UR}_{\Pi,\omega-1}(\mathcal{P})$  produced after eliminating the first block of variables  $X_{[\omega]}$  are bounded respectively by  $O(d)^{k_\omega}$  and  $s^{k_\omega}O(d)^{k_\omega}$ . The number of arithmetic operations in this step is bounded by  $s^{k_\omega}d^{O((k-k_\omega)k_\omega)}$ , and the size of the set  $\text{B}_{\Pi,\omega-1}(\mathcal{P})$  is  $s^{k_\omega+1}d^{O(k_\omega)}$ . Since the cardinality of  $\text{SIGN}_{\Pi,\omega-1}(\mathcal{P})(z)$  is, for every  $z \in \mathbf{R}^{[\omega-1]}$ , bounded by the number of points associated to the univariate representations obtained by substituting  $z$  to the parameters in the elements of  $\text{UR}_{\Pi,\omega-1}(\mathcal{P})$ ,  $\#(\text{SIGN}_{\Pi,\omega-1}(\mathcal{P})(z))$  is  $s^{k_\omega}O(d)^{k_\omega}$ .

An easy inductive argument shows that the number of univariate representations in  $\text{UR}_{\Pi,i}(\mathcal{P})$  produced after eliminating the  $(i+1)$ -th block of variables is bounded by

$$s^{(k_\omega+1)\dots(k_{i+2}+1)k_{i+1}}d^{O(k_\omega)\dots O(k_{i+1})}.$$

By a similar argument, one can show that the degrees of the parametrized univariate representations in  $\text{UR}_{\Pi,i}(\mathcal{P})$  are bounded by  $d^{O(k_\omega)\dots O(k_{i+1})}$ . The complexity in  $D$  is bounded by

$$s^{(k_\omega+1)\dots(k_{i+1}+1)}d^{(k_1+\dots+k_i+2(\omega-i))O(k_\omega)\dots O(k_{i+1})},$$

since the arithmetic is done in a polynomial ring with  $k_1 + \dots + k_i + 2(\omega - i)$  variables.

A similar inductive argument shows that the the size of the set  $\text{B}_{\Pi,i}(\mathcal{P})$  is bounded by  $s^{(k_\omega+1)\dots(k_{i+1}+1)}d^{O(k_\omega)\dots O(k_{i+1})}$ , and their degrees are bounded by  $d^{O(k_\omega)\dots O(k_i)}$ .

The above analysis shows that the size of the set of pairs  $(\mathcal{T}_u, \mathcal{P}_u)$ , constructed at the end of the Substitution Phase is

$$s^{(k_\omega+1)\dots(k_1+1)}d^{O(k_\omega)\dots O(k_1)},$$

and the degrees are bounded by  $d^{O(k_\omega)\cdots O(k_1)}$ . It should also be clear that the number of arithmetic operations in  $D$  for the Substitution Phase is equally bounded by

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}.$$

Since the number of triangular systems  $\mathcal{T}$  is

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)},$$

and each call to Algorithm 12.19 (Triangular Sign Determination) takes time

$$d^{\omega O(k_\omega)\cdots O(k_1)} = d^{O(k_\omega)\cdots O(k_1)},$$

the time taken for the Sign Determination Phase, is

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}.$$

The time required to construct  $\text{SIGN}_\Pi(\mathcal{P})$  is again bounded by

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}.$$

Thus the total time bound for the elimination and sign determination phase is

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}.$$

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\cdots O(k_1)}$ .  $\square$

*Remark 14.12.* In fact, Algorithm 14.2 (Block Structured Signs) does not only compute  $\text{SIGN}_\Pi(\mathcal{P})$ , it also produces the set  $\mathcal{V}$  specifying a  $\Pi$ -set of sampling points for  $\mathcal{P}$  described at the end of Notation 14.10.  $\square$

We have proved the following result:

**Theorem 14.13.** *Let  $\mathcal{P}$  be a set of at most  $s$  polynomials each of degree at most  $d$  in  $k$  variables with coefficients in a real closed field  $\mathbb{R}$ , and let  $\Pi$  denote a partition of the list of variables  $(X_1, \dots, X_k)$  into blocks  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  has size  $k_i$ ,  $1 \leq i \leq \omega$ . Then the size of the set  $\text{SIGN}_\Pi(\mathcal{P})$  is bounded by*

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}.$$

Moreover, there exists an algorithm which computes this set with complexity

$$s^{(k_\omega+1)\cdots(k_1+1)}d^{O(k_\omega)\cdots O(k_1)}$$

in  $D$ , where  $D$  is the ring generated by the coefficients of  $\mathcal{P}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\cdots O(k_1)}$ .

Using the set  $\text{SIGN}_\Pi(\mathcal{P})$ , it is now easy to solve the general decision problem, which is to design a procedure to decide the truth or falsity of a  $(\mathcal{P}, \Pi)$ -sentence.

*Algorithm 14.3. [General Decision]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$
- **Input:** a  $(\mathcal{P}, \Pi)$ -sentence  $\Phi$ , where  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , and  $\Pi$  is a partition of the variables  $X_1, \dots, X_k$  into blocks  $X_{[1]}, \dots, X_{[\omega]}$ .
- **Output:** 1 if  $\Phi$  is true and 0 otherwise.
- **Complexity:**  $s^{(k_\omega+1)\dots(k_1+1)} d^{O(k_\omega)\dots O(k_1)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$ ,  $d$  is a bound on their degree, and  $k_i$  is the number of elements of  $X_{[i]}$ .
- **Procedure:**
  - Compute  $\text{SIGN}_\Pi(\mathcal{P})$ .
  - Trying all possibilities, decide whether

$$\text{Qu}_1\sigma_1 \in \text{SIGN}_\Pi(\mathcal{P}) \quad \text{Qu}_2\sigma_2 \in \sigma_1 \dots \text{Qu}_\omega\sigma_\omega \in \sigma_{\omega-1} \quad F^*(\sigma_\omega) = \text{True},$$

which is clearly a finite verification.

**Proof of correctness:** Follows from the properties of  $\text{SIGN}_\Pi(\mathcal{P})$ . □

**Complexity analysis:** Given the complexity of Algorithm 14.2 (Block Structured Signs), the complexity for the general decision algorithm is

$$s^{(k_\omega+1)\dots(k_1+1)} d^{O(k_\omega)\dots O(k_1)}$$

in  $D$ . Note that the evaluation of the boolean formulas are not counted in this model of complexity since we count only arithmetic operations in  $D$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\dots O(k_1)}$ . □

Note that the first step of the computation depend only on  $(\mathcal{P}, \Pi)$  and not on  $\Phi$ . As noted before  $\text{SIGN}_\Pi(\mathcal{P})$  allows to decide the truth or falsity of every  $(\mathcal{P}, \Pi)$ -sentence.

We have proved the following result.

**Theorem 14.14. [General Decision]** *Let  $\mathcal{P}$  be a set of at most  $s$  polynomials each of degree at most  $d$  in  $k$  variables with coefficients in a real closed field  $\mathbb{R}$ , and let  $\Pi$  denote a partition of the list of variables  $(X_1, \dots, X_k)$  into blocks  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  has size  $k_i$ ,  $1 \leq i \leq \omega$ . Given a  $(\mathcal{P}, \Pi)$ -sentence  $\Phi$ , there exists an algorithm to decide the truth of  $\Phi$  with complexity*

$$s^{(k_\omega+1)\dots(k_1+1)} d^{O(k_\omega)\dots O(k_1)}$$

*in  $D$ , where  $D$  is the ring generated by the coefficients of  $\mathcal{P}$ . If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\dots O(k_1)}$ .*

## 14.2 Quantifier Elimination

In our Quantifier Elimination Algorithm, we use a parametrized version of Algorithm 12.8 (Multivariate Sign Determination) to solve the following problem.

**Notation 14.15.** *Let  $D$  be a ring contained in a real closed field  $R$ . A parametrized triangular system with parameters  $Y = (Y_1, \dots, Y_\ell)$  and variables  $T_1, \dots, T_\omega$  is a list  $\mathcal{T} = T_1, T_2, \dots, T_\omega$  where*

$$\begin{aligned} T_1(Y) &\in D[Y, T_1] \\ T_2(Y) &\in D[Y, T_1, T_2] \\ &\vdots \\ T_\omega(Y) &\in D[Y, T_1, \dots, T_\omega]. \end{aligned}$$

Given a parametrized triangular system  $\mathcal{T} = T_1, T_2, \dots, T_\omega$  with parameters  $Y = (Y_1, \dots, Y_\ell)$ , a set of polynomials  $\mathcal{P} \subset D[Y, T_1, \dots, T_\omega]$  and a point  $y \in R^\ell$  such that  $\mathcal{T}(y)$  is zero-dimensional, we denote by  $\text{SIGN}(\mathcal{P}(y), \mathcal{T}(y))$  the list of sign conditions satisfied by  $\mathcal{P}(y)$  at the zeros of  $\mathcal{T}(y)$ . We want to compute a quantifier free formula such that  $\Phi(z)$  holds if and only if

$$\text{SIGN}(\mathcal{P}(z), \mathcal{T}(z)) = \text{SIGN}(\mathcal{P}(y), \mathcal{T}(y)). \quad \square$$

*Algorithm 14.4.* **[Inverse Sign Determination]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a parametrized triangular system of polynomials,  $\mathcal{T}$  with parameters  $Y = (Y_1, \dots, Y_\ell)$ ,
  - a point  $y \in R^\ell$ , specified by a Thom encoding, such that  $\mathcal{T}(y)$  is zero-dimensional,
  - a subset  $\mathcal{P} \subset D[Y, T_1, \dots, T_\omega]$ .
- **Output:**
  - a family  $\mathcal{A}(y) \subset D[Y]$ ,
  - a quantifier free  $\mathcal{A}(y)$ -formula  $\Phi(y)(Y)$  such that for any  $z \in R^\ell$ , the formula  $\Phi(y)(z)$  is true if and only if  $\mathcal{T}(z)$  is zero-dimensional and

$$\text{SIGN}(\mathcal{P}(y), \mathcal{T}(y)) = \text{SIGN}(\mathcal{P}(z), \mathcal{T}(z)).$$

- **Complexity:**  $s^{\ell+1}(d'^\omega d)^{O(\ell)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{T}$  and  $\mathcal{P}$ .
- **Procedure:**
  - Use Algorithm 12.19 (Triangular Sign Determination) to compute  $\text{SIGN}(\mathcal{Q}(y), \mathcal{T}(y))$ . Form the list

$$B(\text{SIGN}(\mathcal{Q}(y), \mathcal{T}(y))) \subset \{0, 1, 2\}^{\mathcal{Q}},$$

using Remark 10.69 and its notation.

- Using Algorithm 12.18 (Parametrized Bounded Algebraic Sampling) with input  $\mathcal{T}_1^1 + \dots + \mathcal{T}_k^2$ , output a finite set  $\mathcal{U}$  of parametrized univariate representations.
- For every  $\alpha \in B(\text{SIGN}(\mathcal{Q}(y) \cup \text{Der}(\mathcal{T}(y)), \mathcal{T}(y)))$  and every  $u = (f, g_0, \dots, g_k) \in \mathcal{U}$ , compute the signed subresultant coefficients of  $f$  and  $\mathcal{Q}_u^\alpha$ , using Algorithm 8.21 (Signed subresultant) and place them in a set  $\mathcal{A}(y) \subset \mathbb{D}[Y]$ .
- Using Algorithm 13.1 (Computing realizable sign conditions), output the set  $\text{SIGN}(\mathcal{A}(y))$  of realizable sign conditions on  $\mathcal{A}(y)$  and the subset  $\Sigma(y)$  of  $\text{SIGN}(\mathcal{A}(y))$  of  $\rho$  such that for every  $z$  in the realization of  $\rho$ , the Tarski-queries of  $f(z, T)$  and  $\mathcal{Q}_u^\alpha(z, T)$  give rise to a list of non-empty sign conditions  $\text{SIGN}(\mathcal{P}(z), \mathcal{T}(z))$  that coincides with  $\text{SIGN}(\mathcal{P}(y), \mathcal{T}(y))$ .
- Output  $\mathcal{A}(y)$  and

$$\Phi(y)(Y) = \bigvee_{\sigma \in \Sigma(y)} \bigwedge_{Q \in \mathcal{A}(y)} \text{sign}(Q(Y)) = \sigma(Q).$$

**Proof of correctness:** It follows from the correctness of Algorithm 12.19 (Triangular Sign Determination), Remark 10.69, Algorithm 12.18 (Parametrized Bounded Algebraic Sampling), Algorithm 8.21 (Signed subresultant) and Algorithm 13.1 (Computing realizable sign conditions).  $\square$

**Complexity analysis:** Suppose that the degree of  $f_i$  is bounded by  $d'$  and the degrees of all the polynomials in  $\mathcal{P}$  are bounded by  $d$ , and that the number of polynomials in  $\mathcal{P}$  is  $s$ . Using the complexity of Algorithm 12.19 (Triangular Sign Determination), the number of arithmetic operations in  $\mathbb{D}$  in Step 1 is bounded by  $s d'^{O(\omega)}$ . The number of elements of  $B(\text{SIGN}(\mathcal{Q}(y), \mathcal{T}(y)))$  is bounded by  $s O(d')^\omega d$ , using Remark 10.69. The number of arithmetic operations in  $\mathbb{D}[Y]$  is bounded by  $s d'^{O(\omega)} d^{O(1)}$ . The degree in  $Y$  in the intermediate computations is bounded by  $d'^{O(\omega)} d^{O(1)}$ , using the complexity of Algorithm 12.19 (Triangular Sign Determination). Using the complexity analyses of Algorithms 8.4 (Addition of multivariate polynomials), 8.5 (Multiplication of multivariate polynomials), and 8.6 (Exact division of multivariate polynomials), the number of arithmetic operations in  $\mathbb{D}$  is bounded by  $s(d'^\omega d)^{O(\ell)}$ . The number of elements in  $\mathcal{A}(y)$  is  $s d'^{O(\omega)} d^{O(1)}$ . Using the complexity of Algorithm 13.1 (Computing realizable sign conditions), the final complexity is  $s^{\ell+1} (d'^\omega d)^{O(\ell)}$ .

If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau(d'^\omega d)^{O(\ell)}$ .  $\square$

We now describe our algorithm for the quantifier elimination problem. We make use of Algorithm 14.2 (Block Structured Signs) and Algorithm 14.4 (Inverse Sign Determination).

Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k, Y_1, \dots, Y_\ell]$  be finite and let  $\Pi$  denote a partition of the list of variables  $X = (X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  is of size  $k_i$ ,  $1 \leq i \leq \omega$ ,  $\sum_{1 \leq i \leq \omega} k_i = k$ . We proceed in the same manner as the algorithm for the general decision problem, starting with the set  $\mathcal{P}$  of polynomials and eliminating the blocks of variables to obtain a set of polynomials  $B_\Pi(\mathcal{P})$  in the variables  $Y$ . For a fixed  $y \in \mathbb{R}^\ell$ , the truth or falsity of the formula  $\Phi(y)$  can be decided from the set  $\text{SIGN}_\Pi(\mathcal{P})(y)$ . We next apply Algorithm 13.1 (Sampling) to the set of polynomials  $B_\Pi(\mathcal{P}) \subset D[Y]$ , to obtain points in every semi-algebraically connected component of a realizable sign condition of  $B_\Pi(\mathcal{P})$ . For each sample point  $y$  so obtained, we determine whether or not  $y$  satisfies the given formula using the set  $\text{SIGN}_\Pi(\mathcal{P})(y)$ . If it does, then we use the Inverse Sign Determination Algorithm with the various  $\mathcal{T}_u, \mathcal{P}_{u,\omega}, y$  as inputs to construct a formula  $\Psi_y(Y)$ . The only problem left is that this formula contains the infinitesimal quantities introduced by the general decision procedure. However we can replace each equality, or inequality in  $\Psi_y(Y)$ , by an equivalent larger formula without the infinitesimal quantities by using the ordering amongst the infinitesimal quantities. We output the disjunction of the formulas  $\Psi_y(Y)$  constructed above.

We now give a more formal description of the algorithm and prove the bounds on the time complexity and the size of the output formula.

*Algorithm 14.5. [Quantifier Elimination]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite subset  $\mathcal{P} \subset D[X_1, \dots, X_k, Y_1, \dots, Y_\ell]$  of  $s$  polynomials of degree at most  $d$ , a partition  $\Pi$  of the list of variables  $X = (X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  is of size  $k_i$ ,  $1 \leq i \leq \omega$ , with  $\sum_{1 \leq i \leq \omega} k_i = k$  and a  $(\mathcal{P}, \Pi)$ -formula  $\Phi(Y)$ .
- **Output:** a quantifier free formula  $\Psi(Y)$  equivalent to  $\Phi(Y)$ .
- **Complexity:**  $s^{(k_\omega+1)\dots(k_1+1)(\ell+1)} d^{O(k_\omega)\dots O(k_1)O(\ell)}$ .
- **Procedure:**
  - Block Elimination Phase: Perform the Block Elimination Phase of Algorithm 14.2 (Block Structured Signs) on the set of polynomials  $\mathcal{P}$ , with  $\omega + 1$  blocks of variables  $(Y, X_{[1]}, \dots, X_{[\omega]})$  to obtain the set  $\mathcal{U}$  consisting of triangular systems  $\mathcal{T}_u$  and the set of polynomials  $\mathcal{P}_{u,\omega+1}$ .
  - Formula Building Phase: For every  $u = (u_1, \dots, u_{\omega+1}) \in \mathcal{U}$  and every point  $y$  associated to  $u_1$ , compute  $\text{SIGN}(\mathcal{T}_u(y), \mathcal{P}_{u,\omega}(y))$ , using Algorithm 12.19 (Triangular Sign Determination). Output the set  $\text{SIGN}_\Pi(\mathcal{P})(y)$  from the set  $\{\text{SIGN}(\mathcal{T}_u(y), \mathcal{P}_{u,\omega}(y)) \mid u \in \mathcal{U}\}$ , and hence decide whether the formula  $\Phi(y)$  is true.
  - If  $\Phi(y)$  is true, apply Algorithm 14.4 (Inverse Sign Determination) with

$$\mathcal{T}_u, \mathcal{P}_{u,\omega}, y$$



as inputs to get the formulas  $\Psi_{u,y}(Y)$ . Let  $\Psi_y(Y) = \bigwedge_u \Psi_{u,y}(Y)$ , and let  $\overline{\Psi(Y)} = \bigvee_y \Psi_y(Y)$ , where the disjunction is over all the  $y$  for which  $\Phi(y)$  is true in the previous step.

- Output  $\Psi(Y) := \text{Remo}_{\varepsilon_1, \delta_1, \dots, \varepsilon_{\omega+1}, \delta_{\omega+1}}(\overline{\Psi(Y)})$  (Notation 14.6).

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 14.3 ([General Decision), Algorithm 14.4 (Inverse Sign Determination), and Proposition 14.7. □

**Complexity analysis:** The elimination phase takes at most

$$s^{(k_{\omega+1}) \dots (k_1+1)(\ell+1)} d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$$

arithmetic operations, and the number of sign conditions produced is also bounded by

$$s^{(k_{\omega+1}) \dots (k_1+1)(\ell+1)} d^{O(k_{\omega}) \dots O(k_1)O(\ell)}.$$

The degrees in the variables  $T_1, \dots, T_{\omega}, T_{\omega+1}, \varepsilon_1, \delta_1, \dots, \varepsilon_{\omega+1}, \delta_{\omega+1}$  in the polynomials produced, are all bounded by  $d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$ .

Invoking the bound on the Algorithm 14.4 (Inverse Sign Determination), and the bound on the number of tuples produced in the elimination phase, which is  $s^{(k_{\omega+1}) \dots (k_1+1)\ell} d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$  we see that the formula building phase takes no more than

$$s^{(k_{\omega+1}) \dots (k_1+1)\ell + \ell} d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$$

operations. Since the degrees of the variables  $\varepsilon_{\omega+1}, \delta_{\omega+1}, \dots, \varepsilon_1, \delta_1$ , are all bounded by  $d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$ , each atom is expanded to a formula of size at most  $d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$ .

The bound on the size of the formula is an easy consequence of the bound on the number of tuples produced in the elimination phase, and the bound on the formula size produced by Algorithm 14.4 (Inverse Sign Determination).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_{\omega}) \dots O(k_1)O(\ell)}$ . □

This proves the following result.

**Theorem 14.16. [Quantifier Elimination]** *Let  $\mathcal{P}$  be a set of at most  $s$  polynomials each of degree at most  $d$  in  $k + \ell$  variables with coefficients in a real closed field  $\mathbb{R}$ , and let  $\Pi$  denote a partition of the list of variables  $(X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  has size  $k_i$ , for  $1 \leq i \leq \omega$ . Given  $\Phi(Y)$ , a  $(\mathcal{P}, \Pi)$ -formula, there exists an equivalent quantifier free formula,*

$$\Psi(Y) = \bigvee_{i=1}^I \bigwedge_{j=1}^{J_i} \left( \bigvee_{n=1}^{N_{i,j}} \text{sign}(P_{ijn}(Y)) = \sigma_{ijn} \right),$$

where  $P_{i_j n}(Y)$  are polynomials in the variables  $Y$ ,  $\sigma_{i_j n} \in \{0, 1, -1\}$ ,

$$\begin{aligned} I &\leq s^{(k_\omega+1)\dots(k_1+1)(\ell+1)} d^{O(k_\omega)\dots O(k_1)O(\ell)}, \\ J_i &\leq s^{(k_\omega+1)\dots(k_1+1)} d^{O(k_\omega)\dots O(k_1)}, \\ N_{i_j} &\leq d^{O(k_\omega)\dots O(k_1)}, \end{aligned}$$

and the degrees of the polynomials  $P_{i_j k}(y)$  are bounded by  $d^{O(k_\omega)\dots O(k_1)}$ . Moreover, there is an algorithm to compute  $\Psi(Y)$  with complexity

$$s^{(k_\omega+1)\dots(k_1+1)(\ell+1)} d^{O(k_\omega)\dots O(k_1)O(\ell)}$$

in  $\mathbb{D}$ , denoting by  $\mathbb{D}$  the ring generated by the coefficients of  $\mathcal{P}$ .

If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\dots O(k_1)O(\ell)}$ .

*Remark 14.17.* Note that, for most natural geometric properties that can be expressed by a formula in the language of ordered fields, the number of alternations of quantifiers in the formula is small (say at most five or six) while the number of variables can be arbitrarily big. A typical illustrative example is the formula describing the closure of a semi-algebraic set. In such situations, using Theorem 14.16, the complexity of quantifier elimination is singly exponential in the number of variables.  $\square$

**Exercise 14.1.** Design an algorithm computing the minimum value (maybe  $-\infty$ ) of a polynomial of degree  $d$  defined on  $\mathbb{R}^k$  with complexity  $d^{O(k)}$ . Make precise how this minimum value is described.

### 14.3 Local Quantifier Elimination

In this section we discuss a variant of Algorithm 14.5 (Quantifier Elimination) whose complexity is slightly better. A special feature of this algorithm is that the quantifier-free formula output will not necessarily be a disjunction of sign conditions, but will have a more complicated nested structure reflecting the logical structure of the input formula.

For this purpose, we need a parametrized version of Algorithm 12.20 (Triangular Thom Encoding). This algorithm will be based on Algorithm 14.6 (Parametrized Sign Determination).

*Algorithm 14.6. [Parametrized Sign Determination]*

- **Structure:** an ordered domain  $\mathbb{D}$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a parametrized triangular system  $\mathcal{T}$  with parameters  $(Y_1, \dots, Y_\ell)$ , and variables  $(X_1, \dots, X_k)$  and a finite set  $\mathcal{Q} \subset \mathbb{D}[Y_1, \dots, Y_\ell, X_1, \dots, X_k]$ .
- **Output:**
  - a finite set  $\mathcal{A} \subset \mathbb{D}[Y]$ , with  $Y = (Y_1, \dots, Y_k)$ .

- for every  $\rho \in \text{SIGN}(\mathcal{A})$ , a list  $\text{SIGN}(\mathcal{Q}, \mathcal{T})(\rho)$  of sign conditions on  $\mathcal{Q}$  such that, for every  $y$  in the realization  $\text{Reali}(\rho)$  of  $\rho$ ,  $\text{SIGN}(\mathcal{Q}, \mathcal{T})(\rho)$  is the list of sign conditions realized by  $\mathcal{Q}(y)$  on the zero set  $Z(y)$  of  $\mathcal{T}(y)$ .
- **Complexity:**  $s^{\ell(\ell+1)+1} d'^{O(k\ell)} d^{O(\ell)}$ , where  $s$  is a bound on the number of polynomials in  $\mathcal{Q}$  and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{T}$  and  $\mathcal{Q}$ .
- **Procedure:**
  - Step 1: Perform Algorithm 12.18 (Parametrized Bounded Algebraic Sampling) with input  $\mathcal{T}_1^2 + \dots + \mathcal{T}_k^2$ , for  $\mathcal{T}_i \in \mathcal{T}$  and output  $\mathcal{U}$ .
  - Step 2: Consider for every  $u = (f, g_0, \dots, g_k) \in \mathcal{U}$  and every  $Q \in \mathcal{Q}$  the finite set  $\mathcal{F}_{u,Q}$  containing  $Q_u$  (Notation 13.8) and all the derivatives of  $f$  with respect to  $T$ , and compute

$$\mathcal{D}_{u,Q} = \text{RElim}_T(f, \mathcal{F}_{u,Q}) \subset \text{D}[Y],$$

using Algorithm 11.19 (Restricted Elimination).

- Step 3: Define  $\mathcal{D} = \bigcup_{u \in \mathcal{U}, Q \in \mathcal{Q}} \mathcal{D}_{u,Q}$ . Perform Algorithm 13.1 (Sampling) with input  $\mathcal{D}$ . Denote by  $\mathcal{S}$  the set of sample points output.
- Step 4: For every sample point  $y$ , perform Algorithm 14.4 (Inverse Sign Determination) and output the set  $\mathcal{A}(y) \subset \text{D}[Y]$ , as well as  $\text{SIGN}(\mathcal{Q}(y), \mathcal{T}(y))$  and  $\Phi(y)(Y)$ .
- Step 5: Define  $\mathcal{A} = \mathcal{D} \cup \bigcup_{y \in \mathcal{S}} \mathcal{A}(y)$ . Compute the set of realizable sign conditions on  $\mathcal{A}$  using Algorithm 13.1 (Sampling).
- Step 6: For every  $\rho \in \text{SIGN}(\mathcal{A})$  denote by  $y$  the sample point of  $\text{Reali}(\rho)$ . Define  $\text{SIGN}(\mathcal{Q}, \mathcal{T})(\rho)$  as  $\text{SIGN}(\mathcal{Q}(y), \mathcal{T}(y))$ , computed by Algorithm 12.19 (Triangular Sign Determination).

**Proof of correctness:** Follows from the correctness of Algorithm 12.18 (Parametrized Bounded Algebraic Sampling), Algorithm 11.19 (Restricted Elimination), Algorithm 13.1 (Sampling), Algorithm 14.4 (Inverse Sign Determination), Algorithm 13.1 (Sampling) and Algorithm 12.19 (Triangular Sign Determination). □

**Complexity analysis:** We estimate the complexity in terms of the number of parameters  $\ell$ , the number of variables  $k$ , the number  $s$  of polynomials in  $\mathcal{P}$ , a bound  $d'$  on the degrees of the polynomials in  $\mathcal{T}$  and a bound  $d$  on the degrees of the polynomials in  $\mathcal{P}$ .

- Step 1: Using the complexity analysis of Algorithm 12.18 (Parametrized Bounded Algebraic Sampling), the complexity of this step is  $d'^{O(k)}$  in the ring  $\text{D}[Y]$ . The polynomials output are of degree  $O(d')^k$  in  $T$  and of degrees  $d^{O(k)}$  in  $Y$ . Finally, the complexity is  $d'^{O(k\ell)}$  in the ring  $\text{D}$ . The number of elements of  $\mathcal{U}$  is  $O(d')^k$ .
- Step 2: The complexity of this step is  $s d'^{O(k\ell)} d^{O(\ell)}$ , using the complexity analysis of Algorithm 11.19 (Restricted Elimination). The number of polynomials output is  $s d'^{O(k)} d^{O(1)}$ .

- Step 3: The complexity of this step is  $s^\ell d'^{O(k\ell)} d^{O(1)}$ , using the complexity analysis of Algorithm 13.1 (Sampling). There are  $s^\ell d'^{O(k\ell)} d^{O(\ell)}$  points output.
- Step 4: For each sample point, the complexity is  $s^{\ell+1} d'^{O(k\ell)} d^{O(\ell)}$  using the complexity analysis of Algorithm 14.4 (Inverse Sign Determination). So the complexity of this step is  $s^{2\ell+1} d'^{O(k\ell)} d^{O(\ell)}$ . The number of elements of  $\mathcal{A}(y)$  is bounded by  $s d'^{O(k)} d^{O(1)}$  and the degrees of the elements of  $\mathcal{A}(y)$  are bounded by  $d'^{O(k)} d^{O(1)}$ .
- Step 5: The number of elements in  $\mathcal{A}$  is  $s^{\ell+1} d'^{O(k\ell)} d^{O(\ell)}$ , and the degrees of the elements of  $\mathcal{A}$  are bounded by  $d'^{O(k)} d^{O(1)}$ . The complexity of this step is  $s^{\ell(\ell+1)} d'^{O(k\ell)} d^{O(\ell)}$ , using the complexity analysis of Algorithm 13.1 (Sampling).
- Step 6: For every  $\rho$ , the complexity is  $s d'^{O(k\ell)} d^{O(\ell)}$ . So the complexity of this step is  $s^{\ell(\ell+1)+1} d'^{O(k\ell)} d^{O(\ell)}$  using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

Finally the complexity is  $s^{\ell(\ell+1)+1} d'^{O(k\ell)} d^{O(\ell)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d'^{O(k\ell)} d^{O(\ell)}$ .  $\square$

We now define parametrized triangular Thom encodings.

A **parametrized triangular Thom encoding** of level  $k$  with parameters  $Y = (Y_1, \dots, Y_\ell)$  specified by  $\mathcal{A}, \rho, \mathcal{T}, \sigma$  is

- a finite subset  $\mathcal{A}$  of  $\mathbb{R}[Y]$ ,
- a sign condition  $\rho$  on  $\mathcal{A}$ ,
- a triangular system of polynomials  $\mathcal{T}$ , where  $\mathcal{T}_i \in \mathbb{R}[Y, X_1, \dots, X_i]$ ,
- a sign condition  $\sigma$  on  $\text{Der}(\mathcal{T})$  such that for every  $y \in \text{Reali}(\rho)$ , there is a zero  $z(y)$  of  $\mathcal{T}(y)$  with triangular Thom encoding  $\rho$ .

*Algorithm 14.7. [Parametrized Triangular Thom Encoding]*

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a parametrized triangular system  $\mathcal{T}$  with parameters  $(Y_1, \dots, Y_\ell)$  and variables  $(X_1, \dots, X_k)$ .
- **Output:**
  - a finite set  $\mathcal{A} \subset D[Y]$ , with  $Y = (Y_1, \dots, Y_k)$ ,
  - for every  $\rho \in \text{SIGN}(\mathcal{A})$ , a list of sign conditions on  $\text{Der}(\mathcal{T})$  specifying for every  $y \in \text{Reali}(\rho)$ , the list of triangular Thom encodings of the roots of  $\mathcal{T}(y)$ .
- **Complexity:**  $d'^{O(k\ell)}$  where  $d'$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ .
- **Procedure:** Apply Algorithm 14.6 (Parametrized Sign Determination) to  $\mathcal{T}$  and  $\text{Der}(\mathcal{T})$ .

**Proof of correctness:** Immediate. □

**Complexity analysis:** The complexity is  $d'^{O(k\ell)}$ , using the complexity of Algorithm 14.6 (Parametrized Sign Determination). The number of elements in  $\mathcal{A}$  is  $d'^{O(k\ell)}$ , and the degrees of the elements of  $\mathcal{A}$  are bounded by  $d'^{O(k)}$ . □

We follow the notations introduced in the last two sections.

Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k, Y_1, \dots, Y_\ell]$  be finite and let  $\Pi$  denote a partition of the list of variables  $X = (X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$ , where the block  $X_{[i]}$  is of size  $k_i, 1 \leq i \leq \omega, \sum_{1 \leq i \leq \omega} k_i = k$ .

Recall that (Notations 14.8 and 14.10) for every  $i = \omega - 1, \dots, 0$ , the elements of  $\text{UR}_{\Pi,i}(\mathcal{P})$ , are parametrized univariate representations in the variable  $T_{i+1}$ , contained in  $\text{D}[Y, X_{[1]}, \dots, X_{[i]}, \varepsilon_{i+1}, \delta_{i+1}][T_{i+1}]^{k_{i+1}+2}$ . Let

$$u = (u_0, \dots, u_{\omega-1}) \in \mathcal{U} = \prod_{i=0}^{\omega-1} \text{UR}_{\Pi,i}(\mathcal{P}),$$

with

$$u_{i-1} = (f^{[i]}, g_0^{[i]}, g_1^{[i]}, \dots, g_{k_i}^{[i]}).$$

Also recall that we denote,

$$\begin{aligned} \mathcal{T}_{u,i} &= (f^{[1]}(T_1), f_{u,1}^{[2]}(T_1, T_2), \dots, f_{u,i-1}^{[i]}(T_1, T_2, \dots, T_i)), \\ \mathcal{T}_u &= (f^{[1]}(T_1), f_{u,1}^{[2]}(T_1, T_2), \dots, f_{u,\omega-1}^{[\omega]}(T_1, T_2, \dots, T_\omega)). \end{aligned}$$

We now introduce the following notation which is used in the description of the algorithm below.

**Notation 14.18.** Let  $u = (u_0, \dots, u_{j-1}) \in \mathcal{U}_i = \prod_{j=0}^{i-1} \text{UR}_{\Pi,j}(\mathcal{P})$ . We denote by  $\mathcal{L}_{u,i}$  the set of all possible triangular Thom encodings of roots of  $\mathcal{T}_{u,i}$  as  $y$  vary over  $\mathbb{R}\langle \varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega \rangle^\ell$ . □

*Algorithm 14.8. [Local Quantifier Elimination]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite subset  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k, Y_1, \dots, Y_\ell]$ , a partition  $\Pi$  of the list of variables  $X = (X_1, \dots, X_k)$  into blocks,  $X_{[1]}, \dots, X_{[\omega]}$  and a  $(\mathcal{P}, \Pi)$ -formula  $\Phi(Y)$ .
- **Output:** a quantifier free formula,  $\Psi(Y)$ , equivalent to  $\Phi(Y)$ .
- **Complexity:**  $s^{(k_\omega+1)\dots(k_1+1)} d^{\ell O(k_\omega)\dots O(k_1)}$  where  $s$  is a bound on the number of elements of  $\mathcal{P}$ ,  $d$  is a bound on the degree of elements of  $\mathcal{P}$ , and  $k_i$  is the size of the block  $X_{[i]}$ .
- **Procedure:**
  - Initialize  $\text{B}_{\mathcal{P},i,\omega}(\mathcal{P}) := \mathcal{P}$ .
  - Block Elimination Phase: Compute

$$\text{B}_{\Pi,i}(\mathcal{P}) = \text{BElim}_{X_{[i+1]}}(\text{Bor}_{\Pi,i+1}(\mathcal{P})),$$

for  $1 \leq i \leq \omega - 1$ , applying repeatedly  $\text{BElim}_{X_{[i+1]}}$ , using Algorithm 14.1 (Block Elimination).

Compute  $\text{UR}_{\Pi,i}(\mathcal{P})$ , for every  $i = \omega - 1, \dots, 0$ . The elements of  $\text{UR}_{\Pi,i}(\mathcal{P})$  are parametrized univariate representations in the variable  $T_{i+1}$ , contained in

$$D[Y, X_{[1]}, \dots, X_{[i]}, \varepsilon_{i+1}, \delta_{i+1}][T_{i+1}]^{k_{i+1}+2}.$$

– For every

$$u = (u_0, \dots, u_{\omega-1}) \in \mathcal{U} = \prod_{i=0}^{\omega-1} \text{UR}_{\Pi,i}(\mathcal{P}),$$

with

$$u_{i-1} = (f^{[i]}, g_0^{[i]}, g_1^{[i]}, \dots, g_{k_i}^{[i]}),$$

compute the corresponding triangular system,

$$\mathcal{T}_u = (f^{[1]}(Y, T_1), f_{u,1}^{[2]}(Y, T_1, T_2), \dots, f_{u,\omega-1}^{[\omega]}(Y, T_1, T_2, \dots, T_\omega)).$$

(see Notation 14.10).

For  $i = 0 \dots \omega - 1$  compute the sets  $\mathcal{L}_{u,i}$ , using Algorithm 14.7 (Parametrized Triangular Thom Encoding) with input  $\mathcal{T}_{u,i}$ .

– Let

$$\Phi(Y) = (\text{Qu}_1 X_{[1]}) \dots (\text{Qu}_\omega X_{[\omega]}) F(X, Y)$$

where  $\text{Qu}_i \in \{\forall, \exists\}$ ,  $Y = (Y_1, \dots, Y_\ell)$  and  $F(X, Y)$  is a quantifier free  $\mathcal{P}$ -formula.

For every atom of the form  $\text{sign}(P) = \sigma$ ,  $P \in \mathcal{P}$  occurring in the input formula  $F$ , and for every

$$u = (u_0, \dots, u_{\omega-1}) \in \mathcal{U} = \prod_{i=0}^{\omega-1} \text{UR}_{\Pi,i}(\mathcal{P}),$$

with

$$u_{i-1} = (f^{[i]}, g_0^{[i]}, g_1^{[i]}, \dots, g_{k_i}^{[i]}),$$

and  $\tau \in \mathcal{L}_{u,\omega}$  compute using Algorithm 14.5 (Quantifier Elimination) a quantifier-free formula  $\phi_{u,\tau}$  equivalent to the formula

$$(\exists T_1, \dots, T_\omega) \bigwedge_{i=1}^{\omega} \text{SIGN}(\text{Der}(f_{u,i-1}^{[i]})) = \tau_i \bigwedge \text{SIGN}(P_{u,\omega}) = \sigma.$$

Let  $F_{u,\tau}$  denote the quantifier-free formula obtained by replacing every atom  $\phi$  in  $F$  by the corresponding formula  $\phi_{u,\tau}$ .

Also, for every

$$u = (u_0, \dots, u_{\omega-1}) \in \mathcal{U} = \prod_{i=0}^{\omega-1} \text{UR}_{\Pi,i}(\mathcal{P}),$$

with

$$u_{i-1} = (f^{[i]}, g_0^{[i]}, g_1^{[i]}, \dots, g_{k_i}^{[i]}), \tau \in \mathcal{L}_{u,\omega}$$

and for every  $j, 1 \leq j \leq \omega$ , compute using Algorithm 14.5 (Quantifier Elimination) a quantifier-free formula  $\psi_{u,\tau,j}$  equivalent to the formula,

$$(\exists T_1, \dots, T_j) \bigwedge_{i=1}^j \text{SIGN}(\text{Der}(f_{u,i-1}^{[i]})) = \tau_i.$$

- For  $u \in \mathcal{U}$  and  $\tau \in \mathcal{L}_{u,\omega}$ , let

$$\Phi_{\omega,u,\tau} = F_{u,\tau}.$$

Compute inductively for  $i$  from  $\omega - 1$  to 0, and for every

$$u = (u_0, \dots, u_{i-1}) \in \mathcal{U}_i = \prod_{j=0}^{i-1} \text{UR}_{\Pi,j}(\mathcal{P}),$$

and  $\tau \in \mathcal{L}_{u,i}$ ,

$$\begin{aligned} \Phi_{i,u,\tau} &= \bigwedge_{(v,\rho), \bar{v}=u, \bar{\rho}=\tau} (\psi_{v,\rho,i+1} \wedge \Phi_{i+1,v,\rho}) \text{ if } \text{Qu}_{i+1} = \exists, \\ &= \bigwedge_{(v,\rho), \bar{v}=u, \bar{\rho}=\tau} (\psi_{v,\rho,i+1} \implies \Phi_{i+1,v,\rho}) \text{ if } \text{Qu}_{i+1} = \forall. \end{aligned}$$

Take  $\overline{\Phi(Y)} = \Phi_0$ .

- Output  $\Psi(Y) = \text{Remo}_{\varepsilon_1, \delta_1, \dots, \varepsilon_\omega, \delta_\omega}(\overline{\Phi(Y)})$  (Notation 14.6).

**Complexity analysis:** It follows from the complexity analysis of Algorithm a14.1 (Block Elimination), Algorithm 14.7 (Parametrized Triangular Thom Encoding) and Algorithm 14.5 (Quantifier Elimination) that the complexity is bounded by  $s^{(k_\omega+1)\dots(k_1+1)} d^{\ell O(k_\omega)\dots O(k_1)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k_\omega)\dots O(k_1)}$ .

Note that the only improvement compared to Algorithm 14.5 (Quantifier Elimination) is that the exponent of  $s$  does not depend on the number of free variables  $\ell$ . Note also that the total number of polynomials in  $Y$  appearing in the formula is  $s^{(k_\omega+1)\dots(k_1+1)} d^{\ell O(k_\omega)\dots O(k_1)}$ . Determining which are the realizable sign conditions on these polynomials would cost  $s^{(\ell+1)(k_\omega+1)\dots(k_1+1)} d^{\ell O(k_\omega)\dots O(k_1)}$ , but this computation is not part of the algorithm. □

We now give an application of Algorithm (Local Quantifier Elimination) 14.8 to the closure of a semi-algebraic set.

Let  $S$  be a semi-algebraic set described by a quantifier free  $\mathcal{P}$ -formula  $F(X)$ , where  $\mathcal{P}$  is a finite set of  $s$  polynomials of degree at most  $d$  in  $k$  variables. The closure of  $S$  is described by the following quantified formula  $\Psi(X)$

$$\forall Z \quad \exists Y \quad \|X - Y\|^2 < Z^2 \wedge F(Y).$$

Note that  $\Psi(X)$  is a first-order formula with two blocks of quantifiers, the first with one variable and the second one with  $k$  variables. Denote by  $\mathcal{R}$  the set of polynomials in  $k$  variables obtained after applying twice Algorithm 14.1 (Block Elimination) to the polynomials appearing in the formula describing the closure of  $S$  in order to eliminate  $Z$  and  $Y$ . These polynomials have the property that the closure of  $S$  is the union of semi-algebraically connected components of sets defined by sign conditions over  $\mathcal{R}$ . According to Theorem 14.16 the set  $\mathcal{R}$  has  $s^{2k+1}d^{O(k)}$  polynomials and each of these polynomials has degree at most  $d^{O(k)}$ . The complexity for computing  $\mathcal{R}$  is  $s^{2(k+1)}d^{O(k)}$ . Note that we cannot ensure that the closure of  $S$  is described by polynomials in  $\mathcal{R}$ . However, performing Algorithm 14.8 (Local Quantifier Elimination) gives a quantifier-free description of the closure of  $S$  in time  $s^{2(k+1)}d^{O(k)}$  by  $s^{2k+1}d^{O(k)}$  polynomials of degree at most  $d^{O(k)}$ .

### 14.4 Global Optimization

We describe an algorithm for finding the infimum of a polynomial on a semi-algebraic set as well as a minimizer if there exists one.

*Algorithm 14.9. [Global Optimization]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:** a finite subset  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a  $\mathcal{P}$ -semi-algebraic set  $S$  described by a quantifier free formula  $\Phi(X)$  and  $F \in D[X_1, \dots, X_k]$ .
- **Output:** the infimum  $w$  of  $F$  on  $S$ , and a minimizer, i.e. a point  $x \in S$  such that  $F(x) = w$  if such a point exists.
- **Complexity:**  $s^{2k+1}d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on degree of  $F$  and of the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Let  $Y$  be a new variable and  $G = Y - F \in D[Y, X_1, \dots, X_k]$ . Denote by  $S' \subset R^{k+1}$  the realization of  $\Phi \wedge G = 0$ .
  - Call Algorithm 14.1 (Block Elimination) with block of variables  $X_1, \dots, X_k$  and set of polynomials  $\mathcal{P} \cup \{G\} \subset D[Y, X_1, \dots, X_k]$ . Let  $\mathcal{B} \subset D[Y]$  denote  $\text{BElim}_X(\mathcal{P} \cup \{G\})$ .
  - Call Algorithm 10.19 (Univariate Sample Points) with input  $\mathcal{B}$  and denote by  $\mathcal{C}$  the set of sample points so obtained. Each element of  $\mathcal{C}$  is a Thom encoding  $(h, \sigma)$ .
  - For each  $y = (h, \sigma) \in \mathcal{C}$ , the points associated to  $\text{UR}_X(\mathcal{P} \cup \{G\})(y)$  intersect every semi-algebraically connected component of every realizable sign condition of the set  $\mathcal{P} \cup \{G\}(y)$  in  $R(\varepsilon, \delta)^k$ . Compute the subset  $\mathcal{C}'$  of elements  $y \in \mathcal{C}$  such that the set of  $\hat{\cdot}$ points associated to  $\text{UR}_X(\mathcal{P} \cup \{G\})(y)$  meets the extension of  $S'$  to  $R(\varepsilon, \delta)$  using Algorithm 12.19 (Triangular Sign Determination).



- If there is no root  $y$  of a polynomial in  $\mathcal{B}$  such that for all  $y' \in \mathcal{C}'$ ,  $y' \geq y$  holds, define  $w$  as  $-\infty$ . Otherwise, define  $w$  as the maximum  $y \in \mathcal{C}$  which is a root of a polynomial in  $\mathcal{B}$  and such that for all  $y' \in \mathcal{C}'$ ,  $y' \geq y$  holds.
- If  $w = (h, \sigma) \in \mathcal{C}'$ , pick  $u = (f, g_0, \dots, g_k) \in \text{UR}_X(\mathcal{P} \cup \{\mathcal{G}\})(w)$  with associated point in the extension of  $S'$  to  $\mathbb{R}(\varepsilon, \delta)$ . Replace  $\delta$  and  $\varepsilon$  by appropriately small elements from the field of quotients of  $\mathbb{D}$  using Algorithm 11.20 (Removal of Infinitesimals) with input  $f$ , its derivatives and the  $P_u, P \in \mathcal{P}$  and using Remark 11.27. Then clear denominators to obtain univariate representation with entries in  $\mathbb{D}[T]$ .

**Proof of correctness:** Follows clearly from the correctness of Algorithm 14.1 (Block Elimination). □

**Complexity analysis:** The call to Algorithm 14.1 (Block Elimination) costs  $s^{k+1} d^{O(k)}$ . The call to Algorithm 10.19 (Univariate Sample Points) costs  $s^{2k} d^{O(k)}$  since there are at most  $s^k d^{O(k)}$  polynomials of degree at most  $d^{O(k)}$ . Each call to Algorithm 12.19 (Triangular Sign Determination) costs  $s d^{O(k)}$  and there are  $s^{2k} d^{O(k)}$  such calls. The call to Algorithm 11.20 (Removal of Infinitesimals) costs  $(s + d^{O(k)}) d^{O(k)}$  which is  $s d^{O(k)}$ . The total complexity is thus  $s^{2k+1} d^{O(k)}$ . □

### 14.5 Dimension of Semi-algebraic Sets

Let  $S$  be a semi-algebraic set described by a quantifier free  $\mathcal{P}$ -formula  $\Phi(X)$

$$S = \{x \in \mathbb{R}^k \mid \Phi(x)\}$$

where  $\mathcal{P}$  is a finite set of  $s$  polynomials in  $k$  variables with coefficients in a real closed field  $\mathbb{R}$ . We denote by  $\text{SSIGN}(\mathcal{P})$  the set of **strict realizable sign conditions of  $\mathcal{P}$** , i.e. the realizable sign conditions  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$  such that for every  $P \in \mathcal{P}$ ,  $P \neq 0$ ,  $\sigma(P) \neq 0$ .

**Proposition 14.19.** *The dimension of  $S$  is  $k$  if and only if there exists  $\sigma \in \text{SSIGN}(\mathcal{P})$  such that  $\text{Reali}(\sigma) \subset S$ .*

**Proof:** The dimension of  $S$  is  $k$  if and only if there exists a point  $x \in S$  and  $r > 0$  such that  $B(x, r) \subset S$ . The sign condition satisfied by  $\mathcal{P}$  at such an  $x$  is necessarily strict. In the other direction, if the sign condition  $\sigma$  satisfied by  $\mathcal{P}$  at such an  $x$  is strict,  $\text{Reali}(\sigma)$  is open, and contained in  $S$  since  $S$  is defined by a quantifier free  $\mathcal{P}$ -formula. □

It is reasonable to expect that the dimension of  $S$  is  $\geq j$  if and only if the dimension of  $\pi(S)$  is  $j$ , where  $\pi$  is a linear surjection of  $\mathbb{R}^k$  to  $\mathbb{R}^j$ .

Using results from Chapter 13, we are going to prove that using  $j(k - j) + 1$  well chosen linear surjections is enough. Recall that we have defined in Notation 13.26 a family

$$\mathcal{L}_{k,k-j} = \{V_i \mid 0 \leq i \leq j(k - j)\}.$$

of  $j(k - j) + 1$  vector spaces such that any linear subspace  $T$  of  $\mathbb{R}^k$  of dimension  $k' \geq j$  is such that there exists  $0 \leq i \leq j(k - j)$  such that  $V_i$  and  $T$  span  $\mathbb{R}^k$  (see Corollary 13.28). We denote by  $v_k(x)$  the Vandermonde vector

$$(1, x, \dots, x^{k-1}).$$

and by  $V_\ell$  the vector subspace of  $\mathbb{R}^k$  generated by

$$v_k(\ell), v_k(\ell + 1), \dots, v_k(\ell + k - k' - 1).$$

We also defined in Notation 13.26 a linear bijection  $L_{j,i}$  such that  $L_{j,i}(V_i)$  consists of vectors of  $\mathbb{R}^k$  having their last  $j$  coordinates equal to 0. We denote by  $M_{k',\ell} = (d_{k-k',\ell})^{k'} L_{k',\ell}^{-1}$ , with

$$d_{k-k',\ell} = \det(v_{k-k'}(\ell), \dots, v_{k-k'}(\ell + k - k' - 1)),$$

and remarked that  $M_{k',\ell}$  plays the same role as the inverse of  $L_{k',\ell}$  but is with integer coordinates.

We denote by  $\pi_j$  the canonical projection of  $\mathbb{R}^k$  to  $\mathbb{R}^j$  forgetting the first  $k - j$  coordinates.

**Proposition 14.20.** *Let  $0 \leq j \leq k$ . The dimension of  $S$  is  $\geq j$  if and only if there exists  $0 \leq i \leq j(k - j)$  such that the dimension of  $\pi_j(L_{j,i}(S))$  is  $j$ .*

**Proof:** It is clear that if the dimension of  $\pi_j(L_{j,i}(S))$  is  $j$ , the dimension of  $S$  is  $\geq j$ . In the other direction, if the dimension of  $S$  is  $k' \geq j$ , by Proposition 5.53, there exists a smooth point  $x$  of  $S$  of dimension  $k'$  with tangent space denoted by  $T$ . By Corollary 13.28, there exists  $0 \leq i \leq j(k - j)$ , such that  $V_i$  and  $T$  span  $\mathbb{R}^k$ . Since  $L_{j,i}(V_i)$  consists of vectors of  $\mathbb{R}^k$  having their last  $j$  coordinates equal to 0, and  $L_{j,i}(V_i)$  and  $L_{j,i}(T)$  span  $\mathbb{R}^k$ ,  $\pi_j(L_{j,i}(T))$  is  $\mathbb{R}^j$ . Then the dimension of  $\pi_j(L_{j,i}(S))$  is  $j$ .  $\square$

The idea for computing the dimension is simple: check whether the dimension of  $S$  is  $k$  or  $-1$  (i.e. is empty) using Proposition 14.19. If it is not the case, try  $k - 1$  or 0 or, then  $k - 2$  or 1, etc.

*Algorithm 14.10.* **[Dimension]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite subset  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , and a semi-algebraic set  $S$  described by a quantifier free  $\mathcal{P}$ -formula  $\Phi(X)$ .

- **Output:** the dimension  $k'$  of  $S$ .
- **Complexity:**

$$\begin{cases} s^{(k-k')k'} d^{O(k'(k-k'))} & \text{if } k' \geq k/2 \\ s^{(k-k'+1)(k'+1)} d^{O(k'(k-k'))} & \text{if } k' < k/2. \end{cases}$$

where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on their degree.

- **Procedure:**

- Initialize  $j := 0$ .
- (  $\star$  ) Consider the block structure  $\Pi_{k-j}$  with two blocks of variables:  $X_{j+1}, \dots, X_k$  and  $X_1, \dots, X_j$ .
- For every  $i = 0, \dots, j(k-j)$  let  $\mathcal{P}_{k-j,i} = \mathcal{P}(M_{k-j,i})$ , using Notation 13.26 and

$$S_{k-j,i} = \{x \in \mathbb{R}^k \mid \Phi(M_{k-j,i}(x))\}.$$

- Compute  $\text{SIGN}_{\Pi_{k-j}}(\mathcal{P}_{k-j,i})$  using Algorithm 14.2 (Block Structured Signs).
- Defining  $X_{\leq j} = X_1, \dots, X_j$ , compute

$$\text{SSIGN}(\text{BELim}_{X_{\leq j}}(\mathcal{P}_{k-j,i}))$$

using Algorithm 13.2 (Sampling). Note, using Remark 14.12, that every sample point output by Algorithm 14.2 (Block Structured Signs) is above a sample point for  $\text{BELim}_{X_{\leq j}}(\mathcal{P}_{k-j,i})$  output by Algorithm 13.2 (Sampling).

- Check whether one of the strict sign conditions in

$$\text{SSIGN}(\text{BELim}_{X_{\leq j}}(\mathcal{P}_{k-j,i}))$$

is satisfied at some point of  $\pi_{k-j}(S_{k-j,i})$ .

- If one of the strict sign conditions in

$$\text{SSIGN}(\text{BELim}_{X_{\leq j}}(\mathcal{P}_{k-j,i}))$$

is satisfied at some point of  $\pi_{k-j}(S_{k-j,i})$ , output  $k-j$ .

- Consider the block structure  $\Pi_j$  with two blocks of variables:  $X_{k-j+1}, \dots, X_k$  and  $X_1, \dots, X_{k-j}$ .
- For every  $i = 0, \dots, j(k-j)$  let  $\mathcal{P}_{j,i} = \mathcal{P}(M_{j,i})$ , using Notation 13.30 and

$$S_{j,i} = \{x \in \mathbb{R}^k \mid \Phi(M_{j,i}(x))\}.$$

- Compute  $\text{SIGN}_{\Pi_j}(\mathcal{P}_{j,i})$  using Algorithm 14.2 (Block Structured Signs).
- Defining  $X_{\leq k-j} = X_1, \dots, X_{k-j}$ , compute

$$\text{SSIGN}(\text{BELim}_{X_{\leq k-j}}(\mathcal{P}_{j,i}))$$

using Algorithm 13.2 (Sampling). Note, using Remark 14.12, that every sample point output by Algorithm 14.2 (Block Structured Signs) is above a sample point for  $\text{BELim}_{X \leq k-j}(\mathcal{P}_{j,i})$  output by Algorithm 13.2 (Sampling).

- Check whether one of the strict sign conditions in

$$\text{SSIGN}(\text{BELim}_{X \leq k-j}(\mathcal{P}_{j,i}))$$

is satisfied at some point of  $\pi_j(S_{j,i})$ .

- If for every  $i = 0 \dots j(k-j)$  none of the strict sign conditions in

$$\text{SSIGN}(\text{BELim}_{X \leq k-j}(\mathcal{P}_{j,i}))$$

is satisfied at some point of  $\pi_j(S_{j,i})$ , output  $j-1$ .

- Otherwise define  $j := j+1$  and go to  $(\star)$ .

**Proof of correctness:** Follows clearly from Proposition 14.19, Proposition 14.20, the correctness of of Algorithm 14.1 (Block Elimination), Algorithm 13.2 (Sampling). □

**Complexity analysis:** There are at most  $(k+1)/2$  values of  $j$  considered in the algorithm.

For a given  $j$ , the complexity of the call to Algorithm 14.2 (Block Structured Signs) performed is  $s^{(j+1)(k-j+1)}d^{O(j(k-j))}$ , using the complexity analysis of Algorithm 14.2 (Block Structured Signs).

The call to Algorithm 13.2 (Sampling) for  $\text{BELim}_{X \leq j}(\mathcal{P}_{k-j,i})$ , has complexity  $s^{(j+1)(k-j+1)}d^{O(j(k-j))}$ , using the complexity analysis of Algorithm 14.1 (Block elimination) and 13.2 (Sampling), since the number of polynomials is  $s^{j+1}d^{O(j)}$ , their degrees are  $d^{O(j)}$  and their number of variables is  $k-j$ .

Similarly, the call to Algorithm 13.2 (Sampling) for  $\text{BELim}_{X \leq k-j}(\mathcal{P}_{j,i})$ , has complexity  $s^{(j+1)(k-j+1)}d^{O(j(k-j))}$ , using the complexity analysis of Algorithm 14.1 (Block elimination) and 13.2 (Sampling), since the number of polynomials is  $s^{k-j+1}d^{O(k-j)}$ , their degrees are  $d^{O(k-j)}$  and their number of variables is  $j$ .

Finally the total cost of the algorithm is

$$\begin{cases} s^{(k-k')k'}d^{O(k'(k-k'))} & \text{if } k' \geq k/2 \\ s^{(k-k'+1)(k'+1)}d^{O(k'(k-k'))} & \text{if } k' < k/2. \end{cases}$$

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k'(k-k'))}$ .

Note that this complexity result is output sensitive, which means that the complexity depends on the output of the algorithm. □

## 14.6 Bibliographical Notes

The idea of designing algorithms taking into account the block structure is due to Grigor'ev [76], who achieved doubly exponential complexity in the number of blocks for the general decision problem. It should be noted that for a fixed value of  $\omega$ , this is only singly exponential in the number of variables. Heintz, Roy and Solerno [85] and Renegar [133] extended this result to quantifier elimination. Renegar's [133] algorithms solved the general decision problem in time  $(s d)^{O(k_\omega) \cdots O(k_1)}$ , and the quantifier elimination problem in time  $(s d)^{O(k_\omega) \cdots O(k_1) O(\ell)}$ .

Most of the results presented in this chapter are based on [13]. In terms of algebraic complexity (the degree of the equations), the complexity of quantifier elimination presented here is similar to [133]. However the bounds in this chapter are more precise in terms of combinatorial complexity (the dependence on the number of equations). Similarly, the complexity of Algorithm 14.10, coming from [19] improves slightly the result of [163] which computes the dimension of a semi-algebraic set with complexity  $(s d)^{O(k'(k-k'))}$ .

The local quantifier elimination algorithm is based on results in [12].

---

## Computing Roadmaps and Connected Components of Algebraic Sets

In this chapter, we compute roadmaps and connected components of algebraic sets. Roadmaps provide a way to count connected components and to decide whether two points belong to the same connected component. Done in a parametric way the roadmap algorithm also gives a description of the semi-algebraically connected components of an algebraic set. The complexities of the algorithms given in this chapter are much better than the one provided by cylindrical decomposition in Chapter 11 (single exponential in the number of variables rather than doubly exponential).

We first define roadmaps. Let  $S$  be a semi-algebraic set. As usual, we denote by  $\pi$  the projection on the  $X_1$ -axis and set  $S_x = \{y \in \mathbb{R}^{k-1} \mid (x, y) \in S\}$ .

A **roadmap** for  $S$  is a semi-algebraic set  $M$  of dimension at most one contained in  $S$  which satisfies the following roadmap conditions:

- $\text{RM}_1$  For every semi-algebraically connected component  $D$  of  $S$ ,  $D \cap M$  is semi-algebraically connected.
- $\text{RM}_2$  For every  $x \in \mathbb{R}$  and for every semi-algebraically connected component  $D'$  of  $S_x$ ,  $D' \cap M \neq \emptyset$ .

The construction of roadmaps is based on the critical point method, using properties of pseudo-critical values provided in Section 15.1. In Section 15.2 we give an algorithm constructing a roadmap for  $\text{Zer}(Q, \mathbb{R}^k)$ , for  $Q \in \mathbb{R}[X_1, \dots, X_k]$ . As a consequence, we get an algorithm for computing the number of connected components (the zero-th Betti number) of an algebraic set, with single exponential complexity.

In Section 15.3 we obtain an algorithm giving a semi-algebraic description of the semi-algebraically connected components of an algebraic set. The idea behind the algorithm is simple: we perform parametrically the roadmap algorithm with a varying input point.

### 15.1 Pseudo-critical Values and Connectedness

We consider a semi-algebraic set  $S$  as the collection of its fibers  $S_x, x \in \mathbb{R}$ . In the smooth bounded case, critical values of  $\pi$  are the only places where the number of connected components in the fiber can change.

More precisely, we can generalize Proposition 7.6 to the case of a general real closed field.

**Proposition 15.1.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular bounded algebraic hypersurface,  $[a, b]$  such that  $\pi$  has no critical value in  $[a, b]$ , and  $d \in [a, b]$ .*

- a) *The number of semi-algebraically connected components of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  and  $\text{Zer}(Q, \mathbb{R}^k)_d$  are the same.*
- b) *Let  $S$  be a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . Then, for every  $d \in [a, b]$ ,  $S_d$  is semi-algebraically connected.*

Proposition 15.1 immediately implies.

**Proposition 15.2.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded non-singular algebraic hypersurface and  $[a, b]$  such that  $\pi$  has no critical value in  $[a, b]$ . Let  $S$  be a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . Then, for every  $d \in [a, b]$ ,  $S_d$  is semi-algebraically connected.*

**Proposition 15.3.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a non-singular algebraic hypersurface and  $S$  a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . If  $S_{[a,b]}$  is not semi-algebraically connected then  $b$  is a critical value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof of Proposition 15.1:** Over the reals (the case  $\mathbb{R} = \mathbb{R}$ ), the two properties are true according to Proposition 7.6.

We now prove that Properties a and b hold for a general real closed field, using Theorem 5.46 (Semi-algebraic triviality) and the transfer principle (Theorem 2.80).

We first prove Property a.

Let  $\{m_1, \dots, m_N\}$  be a list of all monomials in the variables  $x_1, \dots, x_k$  with degree at most the degree of  $Q$ . To an element  $\text{cof} = (c_1, \dots, c_N)$  of  $\mathbb{R}^N$ , we associate the polynomial

$$\text{Pol}(\text{cof}) = \sum_{i=1}^N c_i m_i.$$

Denoting by  $\text{cof}_i(Q)$  the coefficient of  $m_i$  in  $Q$  and by

$$\text{cof}(Q) = (\text{cof}_1(Q), \dots, \text{cof}_N(Q)),$$

we have  $Q = \text{Pol}(\text{cof}(Q))$ .

Consider the field  $\mathbb{R}_{\text{alg}}$  of real algebraic numbers and the subset  $W \subset \mathbb{R}_{\text{alg}}^{N+2+k}$  defined by

$$W = \{(\text{cof}, a', b', x_1, \dots, x_k) \mid a' \leq x_1 \leq b', \text{Pol}(\text{cof})(x_1, \dots, x_k) = 0\}.$$

The set  $W$  can be viewed as the family of sets  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}_{\text{alg}}^{N+2+k})_{[a', b']}$ , parametrized by  $(\text{cof}, a', b') \in \mathbb{R}_{\text{alg}}^{N+2}$ . We also consider the subset  $W' \subset \mathbb{R}_{\text{alg}}^{N+1+k}$  defined by

$$W' = \{(\text{cof}, d', x_1, \dots, x_k) \mid \text{Pol}(\text{cof})(d', \dots, x_k) = 0\}.$$

The set  $W'$  can be viewed as the family of sets  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}_{\text{alg}}^{N+1+k})_{d'}$ , parametrized by  $(\text{cof}, d') \in \mathbb{R}_{\text{alg}}^{N+1}$ . According to Theorem 5.46 (Hardt's triviality) applied to  $W$  (resp.  $W'$ ), there is a finite partition  $\mathcal{A}$  (resp.  $\mathcal{B}$ ) of  $\mathbb{R}_{\text{alg}}^{N+2}$  (resp.  $\mathbb{R}_{\text{alg}}^{N+1}$ ) into semi-algebraic sets, and for every  $A \in \mathcal{A}$  (resp.  $B \in \mathcal{B}$ ) the sets  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}_{\text{alg}}^{N+2+k})_{[a', b]}$  (resp.  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}_{\text{alg}}^{N+1+k})_{d'}$ ) are semi-algebraically homeomorphic as  $(\text{cof}, a', b')$  varies in  $A$  (resp.  $(\text{cof}, d')$  varies in  $B$ ). Hence, they have the same number of bounded semi-algebraically connected components  $\ell(A)$  (resp.  $\ell(B)$ ).

Using the transfer principle (Theorem 2.80), for every real closed field  $\mathbb{R}$  and every  $(\text{cof}, a', b') \in \text{Ext}(A, \mathbb{R})$  (resp.  $(\text{cof}, d') \in \text{Ext}(B, \mathbb{R})$ ), the set  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}^{N+2+k})_{[a', b]}$  has  $\ell(A)$  (resp.  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}^{N+1+k})_{d'}$  has  $\ell(B)$ ) bounded semi-algebraically connected components. Moreover, since the connected components of

$$W_A = \{(\text{cof}, a', b', x_1, \dots, x_k) \in W \mid (\text{cof}, a', b') \in A\}$$

are semi-algebraic sets defined over  $\mathbb{R}_{\text{alg}}$ , there exists, for every  $A \in \mathcal{A}$ ,  $\ell(A)$  quantifier free formulas

$$\Phi_1(A)(\text{cof}, a', b', x_1, \dots, x_k), \dots, \Phi_{\ell(A)}(A)(\text{cof}, a', b', x_1, \dots, x_k),$$

such that for every real closed field  $\mathbb{R}$  and for every  $(\text{cof}, a', b') \in \text{Ext}(A, \mathbb{R})$  the semi-algebraic sets

$$C_j = \{(x_1, \dots, x_k) \in \mathbb{R}^k \mid \Phi_j(A)(\text{cof}, a', b', x_1, \dots, x_k)\}$$

for  $1 \leq j \leq \ell(A)$  are the bounded semi-algebraically connected components of  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}^{N+2+k})_{[a', b]}$ .

Let  $A$  (resp.  $B$ ) be the set of the partition  $\mathcal{A}$  (resp.  $\mathcal{B}$ ) such that  $\text{cof}(Q), a, b \in \text{Ext}(A, \mathbb{R})$  (resp.  $(\text{cof}(Q), d) \in \text{Ext}(B, \mathbb{R})$ ), and let  $E$  be the semi-algebraic set of  $(\text{cof}, a', b', d') \in (\mathbb{R}_{\text{alg}})^{N+3}$  such that  $(\text{cof}, a', b') \in A$ ,  $(\text{cof}, d') \in B$ ,  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}_{\text{alg}}^{N+2+k})$  is a non-singular algebraic hypersurface,  $\pi$  has no critical value over  $[a', b']$ , and  $a' < d' < b'$ . Using the transfer principle (Theorem 2.80), the set  $E$  is non-empty since  $\text{Ext}(E, \mathbb{R})$  is non-empty, and hence  $\text{Ext}(E, \mathbb{R})$  is non-empty.

Given  $(\text{cof}, a', b', d') \in \text{Ext}(E, \mathbb{R})$ , the number of bounded connected components of  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}^{N+2+k})_{[a', b]}$  is equal to the number of bounded connected components of  $\text{Zer}(\text{Pol}(\text{cof}), \mathbb{R}^{N+2+k})_{d'}$ , since Property 1 holds for the reals. It follows that  $\ell(A) = \ell(B)$ , so the number of bounded semi-algebraically connected components of  $\text{Zer}(Q, \mathbb{R}^k)_{[a, b]}$  is equal to the number of bounded semi-algebraically connected components of  $\text{Zer}(Q, \mathbb{R}^k)_d$ .



To complete the proof of the proposition, it remains to prove Property b. According to the preceding paragraph, there exist  $j$  such that

$$S = \{(x_1, \dots, x_k) \in \mathbb{R}^k \mid \Phi_j(A)(\text{cof}(Q), a, b, x_1, \dots, x_k)\}.$$

Since Property b is true over the reals, the formula expressing that for every  $(\text{cof}, a', b', d') \in \text{Ext}(E, \mathbb{R})$  the set

$$\{(x_2, \dots, x_k) \in \mathbb{R}^k \mid \Phi_j(A)(\text{cof}(Q), a, b, d', \dots, x_k)\}$$

is non-empty is true over the reals. Using the transfer principle (Theorem 2.80), this formula is thus true over any real closed field. Thus,  $S_d$  is non-empty. □

In the non-smooth case, we again consider  $X_1$ -pseudo-critical values introduced in Chapter 12. These pseudo critical-values will also be the only places where the number of connected components in the fiber can change. More precisely, generalizing Proposition 15.2 and Proposition 15.3, we prove the following two propositions, which play an important role for computing roadmaps.

**Proposition 15.4.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded algebraic set and  $S$  a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . If  $v \in (a, b)$  and  $[a, b] \setminus \{v\}$  contains no  $X_1$ -pseudo-critical value on  $\text{Zer}(Q, \mathbb{R}^k)$ , then  $S_v$  is semi-algebraically connected.*

**Proposition 15.5.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be a bounded algebraic set and let  $S$  be a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ . If  $S_{[a,b]}$  is not semi-algebraically connected, then  $b$  is an  $X_1$ -pseudo-critical value of  $\text{Zer}(Q, \mathbb{R}^k)$ .*

Before proving these two propositions, we need some preparation. Suppose that the polynomial  $Q \in \mathbb{R}[X_1, \dots, X_k]$ , and  $(d_1, \dots, d_k)$  satisfy the following conditions:

- $Q(x) \geq 0$  for every  $x \in \mathbb{R}^k$ ,
- $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$  for some  $c \leq 1, c \in \mathbb{R}$ ,
- $d_1 \geq d_2 \geq \dots \geq d_k$ ,
- $\deg(Q) \leq d_1, \text{tDeg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ .

Let  $\bar{d}_i$  be an even number  $> d_i, i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ .

Let  $G_k(\bar{d}, c) = c^{\bar{d}_1} (X_1^{\bar{d}_1} + \dots + X_k^{\bar{d}_k} + X_2^2 + \dots + X_k^2) - (2k - 1)$ , and note that  $\forall x \in B(0, 1/c) \quad G_k(\bar{d}, c)(x) < 0$ .

Using Notation 12.35, we consider

$$\begin{aligned} \text{Def}(Q, \zeta) &= \zeta G_k(\bar{d}, c) + (1 - \zeta) Q, \\ \text{Def}_+(Q, \zeta) &= \text{Def}(Q, \zeta) + X_{k+1}^2. \end{aligned}$$

The algebraic set  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  has the following property which is not enjoyed by  $\text{Zer}(\text{Def}(Q, \zeta), \mathbb{R}\langle \zeta \rangle^k)$ .

**Lemma 15.6.** *Let  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$  be a bounded algebraic set. For every semi-algebraically connected component  $D$  of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$  there exists a semi-algebraically connected component  $D'$  of  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})_{[a,b]}$  such that  $\lim_\zeta(D') = D \times \{0\}$ .*

**Proof:** Let  $y = (y_1, \dots, y_k)$  be a point of  $\text{Ext}(D, \mathbb{R}\langle \zeta \rangle)$ . Since  $y \in B(0, 1/c)$ , we have  $G_k(\bar{d}, c)(y) < 0$ , hence  $\text{Def}(Q, \zeta)(y) < 0$ . Thus, there exists a unique point  $(y, f(y))$  in  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  for which  $f(y) > 0$  and the mapping  $f$  is semi-algebraically continuous. Moreover for every  $z$  in  $D$ ,  $\text{Def}(Q, \zeta)$  is infinitesimal, and hence  $f(z) \in \mathbb{R}\langle \zeta \rangle$  is infinitesimal over  $\mathbb{R}$ . So,  $\lim_\zeta(z, f(z)) = (z, 0)$ . Fix  $x \in D$  and denote by  $D'$  the semi-algebraically connected component of  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  containing  $(x, f(x))$ . Since  $\lim_\zeta(D')$  is connected (Proposition 12.43), contained in  $\text{Zer}(Q, \mathbb{R}^k)$ , and contains  $x$ , it follows that  $\lim_\zeta(D') \subset D$ . Since  $f$  is semi-algebraic and continuous, and  $D$  is semi-algebraically path connected, for every  $z$  in  $D$ , the point  $(z, f(z))$  belongs to the semi-algebraically connected component  $D'$  of  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  containing  $(x, f(x))$ . Since  $\lim_\zeta(z, f(z)) = (z, 0)$ , we have  $\lim_\zeta(D') = D \times \{0\}$ .  $\square$

**Exercise 15.1.** Prove that for

$$Q = ((X + 1)^2 + Y^2 - 1)((X - 1)^2 + Y^2 - 1)((X - 2)^2 + Y^2 - 4)$$

the statement of Lemma 15.6 is false if  $\text{Def}_+(Q, \zeta)$  is replaced by  $\text{Def}(Q, \zeta)$ .

We are now able to prove Proposition 15.4 and Proposition 15.5.

**Proof of Proposition 15.4:** By Lemma 15.6, there exists  $D'$ , a semi-algebraically connected component of  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})_{[a,b]}$  such that  $D \times \{0\} = \lim_\zeta(D')$ . Since  $[a, b] \setminus \{v\}$  contains no  $X_1$ -pseudo-critical value, there exists an infinitesimal  $\beta$  such that the  $X_1$ -critical values on  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$  in the interval  $[a, b]$ , if they exist, lie in the interval  $[v - \beta, v + \beta]$ .

We claim that  $D'_{[v-\beta, v+\beta]}$  is semi-algebraically connected.

Let  $x, y$  be any two points in  $D'_{[v-\beta, v+\beta]}$ . We show that there exists a semi-algebraic path connecting  $x$  to  $y$  lying within  $D'_{[v-\beta, v+\beta]}$ . Since,  $D'$  itself is semi-algebraically connected, there exists a semi-algebraic path,  $\gamma: [0, 1] \rightarrow D'$ , with  $\gamma(0) = x$ ,  $\gamma(1) = y$ , and  $\gamma(t) \in D', 0 \leq t \leq 1$ . If  $\gamma(t) \in D'_{[v-\beta, v+\beta]}$  for all  $t \in [0, 1]$ , we are done. Otherwise, the semi-algebraic path  $\gamma$  is the union of a finite number of closed connected pieces  $\gamma_i$  lying either in  $D'_{[a, v-\beta]}$ ,  $D'_{[v+\beta, b]}$  or  $D'_{[v-\beta, v+\beta]}$ .

By Proposition 15.2 the connected components of  $D'_{v-\beta}$  (resp.  $D'_{v+\beta}$ ) are in 1-1 correspondence with the connected components of  $D'_{[a,v-\beta]}$  (resp.  $D'_{[v+\beta,b]}$ ) containing them. Thus, we can replace each of the  $\gamma_i$  lying in  $D'_{[a,v-\beta]}$  (resp.  $D'_{[v+\beta,b]}$ ) with endpoints in  $D'_{v-\beta}$  (resp.  $D'_{v+\beta}$ ) by another segment with the same endpoints but lying completely in  $D'_{v-\beta}$  (resp.  $D'_{v+\beta}$ ). We thus obtain a new semi-algebraic path  $\gamma'$  connecting  $x$  to  $y$  and lying inside  $D'_{[v-\beta,v+\beta]}$ .

It is clear that  $\lim_{\zeta} (D'_{[v-\beta,v+\beta]})$  coincides with  $D_v$ . Since  $D'_{[v-\beta,v+\beta]}$  is bounded,  $D_v$  is semi-algebraically connected by Proposition 12.43.  $\square$

**Proof of Proposition 15.5:** By Lemma 15.6, there exists  $D'$ , a semi-algebraically connected component of  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})_{[a,b]}$  such that  $D \times \{0\} = \lim_{\zeta} (D')$ . According to Theorem 5.46 (Hardt's triviality), there exists  $a' \in [a, b)$  such that for every  $d \in [a', b)$ ,  $D'_{[a,d]}$  is not semi-algebraically connected. Hence, by Proposition 12.43,  $D'_{[a,c]}$  is also not semi-algebraically connected for every  $c \in \mathbb{R}\langle \zeta \rangle$  with  $\lim_{\zeta} (c) = d$ . Since  $D'$  is semi-algebraically connected, according to Proposition 15.3, there is an  $X_1$ -critical value  $c$  on  $\text{Zer}(\text{Def}_+(Q, \zeta), \mathbb{R}\langle \zeta \rangle^{k+1})$ , infinitesimally close to  $b$ . Hence  $b$  is an  $X_1$ -pseudo-critical value on  $\text{Zer}(Q, \mathbb{R}^k)$ .  $\square$

## 15.2 Roadmap of an Algebraic Set

We describe the construction of a roadmap  $M$  for a bounded algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  which contains a finite set of points  $\mathcal{N}$  of  $\text{Zer}(Q, \mathbb{R}^k)$ . A precise description of how the construction can be performed algorithmically will follow.

We first construct  $X_2$ -pseudo-critical points on  $\text{Zer}(Q, \mathbb{R}^k)$  in a parametric way along the  $X_1$ -axis. This results in curve segments and their endpoints on  $\text{Zer}(Q, \mathbb{R}^k)$ . The curve segments are continuous semi-algebraic curves parametrized by open intervals on the  $X_1$ -axis, and their endpoints are points of  $\text{Zer}(Q, \mathbb{R}^k)$  above the corresponding endpoints of the open intervals. Since these curves and their endpoints include, for every  $x \in \mathbb{R}$ , the  $X_2$ -pseudo-critical points of  $\text{Zer}(Q, \mathbb{R}^k)_x$ , they meet every connected component of  $\text{Zer}(Q, \mathbb{R}^k)_x$ . Thus the set of curve segments and their endpoints already satisfy  $\text{RM}_2$ . However, it is clear that this set might not be semi-algebraically connected in a semi-algebraically connected component, so  $\text{RM}_1$  might not be satisfied (see Figure 15). We add additional curve segments to ensure that  $\text{Mea}$  is connected by recursing in certain distinguished hyperplanes defined by  $X_1 = z$  for distinguished values  $z$ .

The set of **distinguished values** is the union of the  $X_1$ -pseudo-critical values, the first coordinates of the input points  $\mathcal{N}$  and the first coordinates of the endpoints of the curve segments. A **distinguished hyperplane** is an hyperplane defined by  $X_1 = v$ , where  $v$  is a distinguished value. The input points, the endpoints of the curve segments and the intersections of the curve segments with the distinguished hyperplanes define the set of **distinguished points** .

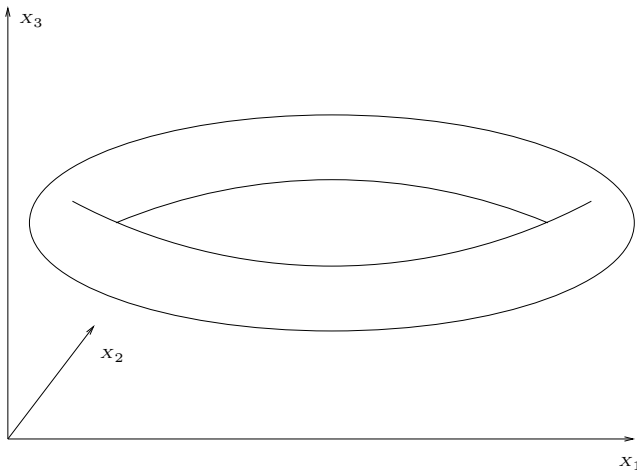
So we have constructed the distinguished values  $v_1 < \dots < v_\ell$  of  $X_1$  among which are the  $X_1$ -pseudo-critical values. Above each interval  $(v_i, v_{i+1})$ , we have constructed a collection of curve segments  $\mathcal{C}_i$  meeting every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_v$  for every  $v \in (v_i, v_{i+1})$ . Above each distinguished value  $v_i$ , we have constructed a set of distinguished points  $\mathcal{N}_i$ . Each curve segment in  $\mathcal{C}_i$  has an endpoint in  $\mathcal{N}_i$  and another in  $\mathcal{N}_{i+1}$ . Moreover, the union of the  $\mathcal{N}_i$  contains  $\mathcal{N}$ .

We then repeat this construction in each distinguished hyperplane  $H_i$  defined by  $X_1 = v_i$  with input  $Q(v_i, X_2, \dots, X_k)$  and the distinguished points in  $\mathcal{N}_i$ .

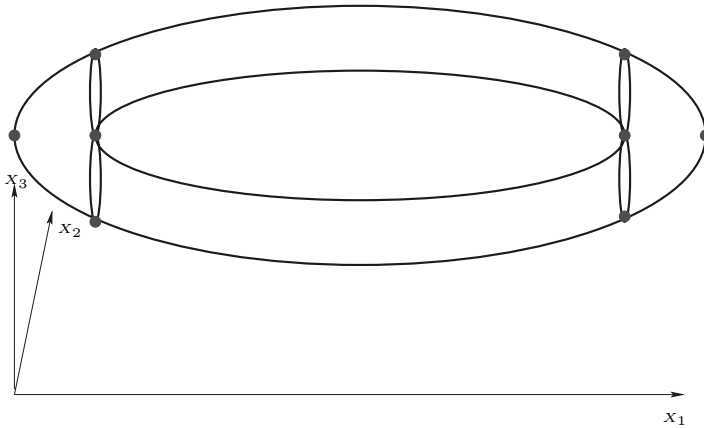
The process is iterated until for

$$I = (i_1, \dots, i_{k-2}), 1 \leq i_1 \leq \ell, \dots, 1 \leq i_{k-2} \leq \ell(i_1, \dots, i_{k-3}),$$

we have distinguished values  $v_{I,1} < \dots < v_{I,\ell(I)}$  along the  $X_{k-1}$  axis with corresponding sets of curve segments and sets of distinguished points with the required incidences between them.



**Fig. 15.1.** A torus in  $\mathbb{R}^3$



**Fig. 15.2.** The roadmap of the torus

**Proposition 15.7.** *The semi-algebraic set  $M$  obtained by this construction is a roadmap for  $\text{Zer}(Q, \mathbb{R}^k)$ .*

The proof of Proposition 15.7 uses the following lemmas.

**Lemma 15.8.** *If  $v \in (a, b)$  is a distinguished value such that  $[a, b] \setminus \{v\}$  contains no distinguished value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$  and  $D$  is a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{[a,b]}$ , then  $M \cap D$  is semi-algebraically connected.*

**Proof:** Since  $[a, b] \setminus \{v\}$  contains no pseudo-critical value of the algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$ , we know, by Proposition 15.4, that  $D_v$  is semi-algebraically connected. Moreover, the points of  $M \cap D$  are connected through curve segments to the points of  $\mathcal{N}_v$ . By induction hypothesis, the points of  $\mathcal{N}_v$  are in the same semi-algebraically connected component of  $D_v$ , since  $D_v$  is semi-algebraically connected.

The construction makes a recursive call at every distinguished hyperplane and hence at  $H_v$ . The input to the recursive call is the algebraic set  $\text{Zer}(Q, \mathbb{R}^k)_v$  and the set of all distinguished points in  $H_v$  which includes the endpoints of the curves in  $M \cap D \cap H_v$ . Hence, by the induction hypothesis they are connected by the roadmap in the slice.

Therefore,  $M \cap D$  is semi-algebraically connected.  $\square$

**Lemma 15.9.** *If  $D$  is a semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ , then  $M \cap D$  is semi-algebraically connected.*

**Proof:** Let  $x, y$  be two points of  $M \cap D$ , and let  $\gamma$  be a semi-algebraic path in  $D$  from  $x$  to  $y$  such that  $\gamma(0) = x$ ,  $\gamma(1) = y$ . We are going to construct another semi-algebraic path from  $x$  to  $y$  inside  $M$ . Let  $\{v_1 < \dots < v_\ell\}$  be the set of distinguished values and choose  $u_i$  such that

$$u_1 < v_1 < u_2 < v_2 < \dots < u_\ell < v_\ell < u_{\ell+1}.$$

There exist a finite number of points of  $\gamma$ ,  $x = x_0, x_1, \dots, x_{N+1} = y$ , with  $\pi(x_i) = u_{n(i)}$ , and semi-algebraic paths  $\gamma_i$  from  $x_i$  to  $x_{i+1}$  such that:

- $\gamma = \bigcup_{0 \leq i \leq N} \gamma_i$ ,
- $\gamma_i \subset D_{[u_{n(i)}, u_{n(i)+1}]}$  or  $\gamma_i \subset D_{[u_{n(i)-1}, u_{n(i)}]}$ .

Let  $D_i$  be the semi-algebraically connected component of  $D_{[u_{n(i)}, u_{n(i)+1}]}$  (resp.  $D_{[u_{n(i)-1}, u_{n(i)}]}$ ) containing  $\gamma_i$ . Since  $D_{i-1} \cap D_i$  is a finite union of semi-algebraically connected components of  $D_{\pi(x_i)}$ ,  $M \cap D_{i-1} \cap D_i$  is not empty. Choose  $y_0 = x, \dots, y_i \in M \cap D_{i-1} \cap D_i, \dots, y_{N+1} = y$ . Then  $y_i$  and  $y_{i+1}$  are in the same semi-algebraically connected component of  $M \cap D$  by Lemma 15.8.  $\square$

**Proof of Proposition 15.7:** We have already seen that  $M$  satisfies  $RM_2$ . We now prove that  $M$  satisfies  $RM_1$ .

The proof is by induction on the dimension of the ambient space. In the case of dimension one, the roadmap properties are obviously true for the set we have constructed. Now assume that the construction gives a roadmap for all dimensions less than  $k$ . That the construction gives a roadmap for dimension  $k$  follows from the following two lemmas. Lemma 15.8 and Lemma 15.9.  $\square$

We now describe precisely a way of performing algorithmically the preceding construction.

In our inductive construction of the roadmap, we are going to use the following specification describing points and curve segments:

A **real univariate triangular representation**  $\mathcal{T}, \sigma, u$  of level  $i - 1$  consists of:

- a triangular Thom encoding  $\mathcal{T}, \sigma$  specifying  $(z, t) \in \mathbb{R}^i$  with  $z \in \mathbb{R}^{i-1}$
- a parametrized univariate representation

$$u(X_{<i}) = (\mathcal{T}_i(X_{<i}, T), g_0(X_{<i}, T), g_i(X_{<i}, T), \dots, g_k(X_{<i}, T)),$$

with parameters  $X_{<i} = (X_1, \dots, X_{i-1})$  (see Definition page 481).

The point associated to  $\mathcal{T}, \sigma, u$  is

$$\left( z, \frac{g_i(z, t)}{g_0(z, t)}, \dots, \frac{g_k(z, t)}{g_0(z, t)} \right).$$

A real univariate triangular representation  $\mathcal{T}, \sigma, u$  is **above** the triangular Thom encoding  $\mathcal{T}', \sigma'$  if  $\mathcal{T}' = \mathcal{T}_1, \dots, \mathcal{T}_{i-1}, \sigma' = \sigma_1, \dots, \sigma_{i-1}$ .

It will be useful to compute the  $i$ -th projection of a point specified by a real univariate representation.

*Algorithm 15.1. [Projection]*

- **Structure:** a domain  $D$  contained in a field  $K$ .

• **Input:**

a real univariate triangular representation  $\mathcal{T}, \sigma, u$  of level  $i - 1$  with coefficients in  $D$ . We denote by  $z$  the root of  $\mathcal{T}$  specified by  $\sigma$  and by  $x$  the point associated to  $\mathcal{T}, \sigma, u$ .

• **Output:** a Thom encoding  $\text{proj}_i(u), \text{proj}_i(\tau)$  specifying the projection of the point associated to  $\mathcal{T}, \sigma, u$  on the  $X_i$  axis.

• **Complexity:**  $d^{O(i)}$ , where  $d$  is a bound on the degree of the univariate representation and a bound on the degrees of the polynomials in  $\mathcal{T}$ .

• **Procedure:**

- Compute the resultant  $\text{proj}_i(u)$  of  $\mathcal{T}_i(X_{<i}, T)$ , and

$$X_i g_0(X_{<i}, T) - g_i(X_{<i}, T)$$

with respect to  $T$ , using Algorithm 8.21 (Signed subresultant).

- Compute the Thom encoding of the root of  $\text{proj}_i(u)$  which is the  $i$ -th coordinate of  $x$  as follows: let  $d$  be the smallest even number not less than the degree of  $\text{proj}_i(u)$  with respect to  $X_i$ , and compute the sign of the derivatives of

$$g_0(X_{<i}, T)^d \text{proj}_i(u) \left( \frac{g_i(X_{<i}, T)}{g_0(X_{<i}, T)} \right)$$

with respect to  $T$  at the root  $z$  of  $\mathcal{T}$  specified by  $\sigma$ . This gives the Thom encoding  $\text{proj}_i(\tau)$  of the  $i$ -th coordinate of  $x$ . This is done using Algorithm 12.19 (Triangular Sign Determination).

**Proof of correctness:** Immediate. □

**Complexity analysis:** The complexity is  $d^{O(ki)}$  using the complexity of Algorithm 12.19 (Triangular Sign Determination).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(i)}$ . □

Let  $\mathcal{V}_1, \tau_1, \mathcal{V}_2, \tau_2$  be two triangular Thom encodings above  $\mathcal{T}, \sigma$ . We denote by  $z = (z_1, \dots, z_{i-1}) \in \mathbb{R}^{i-1}$  the point specified by  $\mathcal{T}, \sigma$  and by  $(z, a), (z, b)$  the points specified by  $\mathcal{V}_1, \tau_1$  and  $\mathcal{V}_2, \tau_2$  (see Definition page 496).

A **curve segment representation**  $u, \rho$  above  $\mathcal{V}_1, \tau_1, \mathcal{V}_2, \tau_2$  is:

- a parametrized univariate representation with parameters  $(X_{\leq i})$

$$u = (f(X_{\leq i}, T), g_0(X_{\leq i}, T), g_{i+1}(X_{\leq i}, T), \dots, g_k(X_{\leq i}, T)),$$

- a sign condition  $\rho$  on  $\text{Der}(f)$  such that for every  $v \in (a, b)$  there exists a real root  $t(v)$  of  $f(z, v, T)$  with Thom encoding  $\sigma, \rho$  and  $g_0(z, v, t(v)) \neq 0$ .

The **curve segment associated to**  $u, \rho$  is the semi-algebraic function  $h$  which maps a point  $v$  of  $(a, b)$  to the point of  $\mathbb{R}^k$  defined by

$$h(v) = \left( z, v, \frac{g_{i+1}(z, v, t(v))}{g_0(z, v, t(v))}, \dots, \frac{g_k(z, v, t(v))}{g_0(z, v, t(v))} \right).$$

It is a continuous injective semi-algebraic function.

The Curve Segments Algorithm will be the basic building block in our algorithm.

*Algorithm 15.2. [Curve Segments]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a triangular Thom encoding  $\mathcal{T}, \sigma$  with coefficients in  $D$  specifying  $z \in \mathbb{R}^{i-1}$ ,
  - a polynomial  $Q \in D[X_1, \dots, X_k]$ , for which  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a finite set  $\mathcal{N}$  of real univariate triangular representation above  $\mathcal{T}, \sigma$  with coefficients in  $D$  and associated points contained in  $\text{Zer}(Q, \mathbb{R}^k)$ .
- **Output:**
  - an ordered list of triangular Thom encodings  $\mathcal{V}_1, \tau_1, \dots, \mathcal{V}_\ell, \tau_\ell$  above  $\mathcal{T}, \sigma$  specifying points  $(z, v_1), \dots, (z, v_\ell)$  with  $v_1 < \dots < v_\ell$ . The  $v_j$  are called distinguished values.
  - For every  $j = 1, \dots, \ell$ ,
    - a finite set  $\mathcal{D}_j$  of real univariate triangular representations representation above  $\mathcal{V}_j, \tau_j$ . The associated points are called distinguished points.
    - a finite set  $\mathcal{C}_j$  of curve segment representations above  $\mathcal{V}_j, \tau_j, \mathcal{V}_{j+1}, \tau_{j+1}$ . The associated curve segments are called distinguished curves.
    - a list of pairs of elements of  $\mathcal{C}_j$  and  $\mathcal{D}_j$  (resp.  $\mathcal{C}_{j+1}$  and  $\mathcal{D}_j$ ) describing the adjacency relations between distinguished curves and distinguished points.

The distinguished curves and points are contained in  $\text{Zer}(Q, \mathbb{R}^k)_z$ . Among the distinguished values are the first coordinates of the points in  $\mathcal{N}$  as well as the pseudo-critical values of  $\text{Zer}(Q, \mathbb{R}^k)_z$ . The sets of distinguished values, distinguished curves, and distinguished points satisfy the following properties.

- CS<sub>1</sub>: For every  $v \in \mathbb{R}$ , the set of distinguished curve and distinguished points output intersect every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)_{(z, v)}$ .
- CS<sub>2</sub>: For each distinguished curve output over an interval with endpoint a given distinguished value, there exists a distinguished point over this distinguished value which belongs to the closure of the curve segment.



- **Complexity:**  $d^{O(ik)}$ , where  $d$  is a bound on the degree of  $Q$  and  $O(d)^k$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ , the degrees of the univariate representations in  $\mathcal{N}$ , and the number of these univariate representations.

- **Procedure:**

- Step 1: Perform Algorithm 12.10 (Parametrized Multiplication Table) with input  $\overline{\text{Cr}}(Q^2, \zeta)$ , (using Notation 12.46) and parameter  $X_{\leq i}$ . Perform Algorithm 12.15 (Parametrized Limit of Bounded Points) and output  $\mathcal{U}$ .

Consider for every  $u = (f, g_0, g_{i+1}, \dots, g_k) \in \mathcal{U}$  the finite set  $\mathcal{F}_u$  containing  $Q_u$  (Notation 13.8) and all the derivatives of  $f$  with respect to  $T$ , and compute

$$\mathcal{D}_u = \text{RElim}_T(f, \mathcal{F}_u) \subset D[X_{\leq i}],$$

using Algorithm 11.19 (Restricted Elimination). Define  $\mathcal{D} = \bigcup_{u \in \mathcal{U}} \mathcal{D}_u$ .

- Step 2: For every  $T', \tau, u \in \mathcal{N}$ , compute  $\text{proj}_i(u), \text{proj}_i(\tau)$  using Algorithm 15.1 (Projection), add to  $\mathcal{D}$  the polynomial  $\text{proj}_i(u)$ .
- Step 3: Compute the Thom encodings of the zeroes of  $A$ ,  $A \in \mathcal{D}$  above  $\mathcal{T}$ ,  $\sigma$  using Algorithms 12.20 (Triangular Thom Encoding), output their ordered list  $A_1, \alpha_1, \dots, A_\ell, \alpha_\ell$  and the corresponding ordered list  $v_1 < \dots < v_\ell$  of distinguished values using Algorithm 12.21 (Triangular Comparison of Roots). Define  $\mathcal{V}_i, \tau_i = \mathcal{T}, A_i, \sigma, \alpha_i$ .
- Step 4: For every  $j = 1, \dots, \ell$  and every  $(f, g_0, g_i, \dots, g_k), \tau \in \mathcal{N}$  such that  $\text{proj}_i(\tau) = \alpha_j$ , append  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{D}_j$ , using Algorithm 12.19 (Triangular Sign Determination).
- Step 5: For every  $j = 1, \dots, \ell$  and every

$$u = (f, g_0, g_{i+1}, \dots, g_k) \in \mathcal{U},$$

compute the Thom encodings  $\tau$  of the roots of  $f$  above  $\mathcal{T}$ ,  $\sigma$  such that  $\text{proj}_i(\tau) = \alpha_j$ , using Algorithm 12.20 (Triangular Thom Encoding). Append all pairs  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{D}_j$  when the corresponding associated point belongs to  $\text{Zer}(Q, \mathbb{R}^k)_z$ .

- Step 6: For every  $j = 1, \dots, \ell - 1$  and every

$$u = (f, g_0, g_{i+1}, \dots, g_k) \in \mathcal{U},$$

compute the Thom encodings  $\rho$  of the roots of  $f(z, v, T)$  over  $(v_j, v_{j+1})$  using Algorithm 12.22 (Triangular Intermediate Points) and Algorithm 12.20 (Triangular Thom Encoding) and append pairs  $u, \rho$  to  $\mathcal{C}_j$  when the corresponding associated curve is included in  $\text{Zer}(Q, \mathbb{R}^k)_z$ .

- Step 7: Determine adjacencies between curve segments and points. For every point of  $\mathcal{D}_j$  specified by

$$v' = (p, q_0, q_{i+1}, \dots, q_k), \tau', \text{ with } \{p, q_0, q_{i+1}, \dots, q_k\} \subset D[X_{\leq i}][T]$$

and every curve segment representation of  $\mathcal{C}_j$  specified by

$$v = (f, g_0, g_{i+1}, \dots, g_k), \tau, \{f, g_0, g_{i+1}, \dots, g_k\} \subset \mathbb{D}[X_{\leq i}][T],$$

decide whether the associated point  $t$  is adjacent to the associated curve segment as follows: compute the first  $\nu$  such that  $(\partial^\nu g_0 / \partial X_i^\nu)(v_j, t)$  is not zero and decide whether for every  $\ell = i + 1, \dots, k$

$$\frac{\partial^\nu g_\ell}{\partial X_i^\nu}(v_j, t) q_0(t) - \frac{\partial^\nu g_0}{\partial X_i^\nu}(v_j, t) q_\ell(t)$$

is zero. This is done using Algorithm 12.19 (Triangular Sign Determination) above  $\mathcal{T}, \sigma$ .

Repeat the same process for every element of  $\mathcal{D}_{j+1}$  and every curve segment representation of  $\mathcal{C}_j$ .

**Proof of correctness:** It follows from Proposition 12.42, the correctness of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination), Algorithm 15.1, Algorithm 12.22 (Triangular Intermediate Points), Algorithm 12.20 (Triangular Thom Encoding), Algorithm 12.21 (Triangular Comparison of Roots) and Algorithm 12.19 (Triangular Sign Determination).  $\square$

### Complexity analysis:

- Step 1: This step requires  $d^{O(i(k-i))}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination). There are  $d^{O(k-i)}$  parametrized univariate representations computed in this step and each polynomial in these representations has degree  $O(d)^{k-i}$ .
- Step 2: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 15.1 (Projection).
- Step 3: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.20 (Triangular Thom Encoding).
- Step 4: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).
- Step 5: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.20 (Triangular Thom Encoding).
- Step 6: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.22 (Triangular Intermediate Points), Algorithm 12.20 (Triangular Thom Encoding).
- Step 7: This step requires  $d^{O(ik)}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

Thus, the complexity is  $d^{O(ik)}$ . The number of distinguished values is bounded by  $d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(ik)}$ .  $\square$

Given a polynomial  $Q$  and a set of real univariate representations  $\mathcal{N}$ , we denote by  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \mathcal{N})$  a roadmap of  $\text{Zer}(Q, \mathbb{R}^k)$  which contains the points associated to  $\mathcal{N}$ .

We now describe a recursive roadmap algorithm for bounded algebraic sets.

*Algorithm 15.3. [Bounded Algebraic Roadmap]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a triangular Thom encoding  $\mathcal{T}, \sigma$  with coefficients in  $D$  specifying  $z \in \mathbb{R}^i$ ,
  - a polynomial  $Q \in D[X_1, \dots, X_k]$ , for which  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a finite set  $\mathcal{N}$  of real univariate representation  $u, \tau$  above  $\mathcal{T}, \sigma$  with coefficients in  $D$  with associated points contained in  $\text{Zer}(Q, \mathbb{R}^k)_z$ .
- **Output:** a roadmap  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z, \mathcal{N})$  which contains the points associated to  $\mathcal{N}$ .
- **Complexity:**  $d^{O(k^2)}$ , where  $d$  is a bound on the degree of  $Q$  and  $O(d)^k$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ , the degrees of the univariate representations in  $\mathcal{N}$ , and the number of these univariate representations.
- **Procedure:**
  - Call Algorithm 15.2 (Curve Segments), output  $\ell$  and, for every  $j = 1, \dots, \ell$ ,  $A_j, \alpha_j, \mathcal{D}_j$  and  $\mathcal{C}_j$ .
  - For every  $j = 1, \dots, \ell$ , call Algorithm 15.3 (Bounded Algebraic Roadmap) recursively, with input  $\mathcal{T}, A_j, \sigma, \alpha_j$ , specifying  $(z, v_j)$ ,  $Q$  and  $\mathcal{D}_j$ .

**Proof of correctness:** The correctness of the algorithm follows from Proposition 15.7 and the correctness of Algorithm 15.2 (Curve Segments).  $\square$

**Complexity analysis:** In the recursive calls to Algorithm 15.3 (Bounded Algebraic Roadmap), the number of triangular systems considered is at most  $d^{O(k^2)}$  and the triangular systems involved have polynomials of degree  $O(d)^k$ . Thus the total number of arithmetic operations in  $D$  is bounded by  $d^{O(k^2)}$  using the complexity analysis of Algorithm 15.2 (Curve Segments).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .  $\square$

Since  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z, \mathcal{N})$  contains  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z)$ , it is possible to extract from  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z, \{u, \tau\})$  a path connecting the point  $p$  associated to  $u, \tau$  to  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z)$ .

*Algorithm 15.4. [Bounded Algebraic Connecting]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a triangular Thom encoding  $\mathcal{T}, \sigma$  with coefficients in  $D$  specifying  $z \in \mathbb{R}^i$ ,
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  for which  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a real univariate triangular representation  $\mathcal{V}, \tau, u$  above  $\mathcal{T}, \sigma$  with coefficients in  $D$ , with associated point  $p$  contained in  $\text{Zer}(Q, \mathbb{R}^k)_z$ .
- **Output:** a path  $\gamma(p) \subset \text{Zer}(Q, \mathbb{R}^k)_z$  connecting  $p$  to a distinguished point of  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k)_z)$ .
- **Complexity:**  $d^{O(k^2)}$ , where  $d$  is a bound on the degree of  $Q$  and  $O(d)^k$  is a bound on the degrees of the polynomials in  $\mathcal{T}$  and the degree of the real univariate triangular representation  $\mathcal{V}, \tau, u$ .
- **Procedure:** Call Algorithm 15.3 (Bounded Algebraic Roadmap) with input  $Q, \mathcal{T}, \sigma$  and  $\{\mathcal{V}, \tau, u\}$ , and extract  $\gamma(p)$  from  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \{\mathcal{V}, \tau, u\})$ .

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 15.3 (Bounded Algebraic Roadmap). □

**Complexity analysis:**The total number of arithmetic operations in  $D$  is bounded by  $d^{O(k^2)}$ , using the complexity analysis of Algorithm 15.3 (Bounded Algebraic Roadmap).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

*Remark 15.10.* Note that the connecting path  $\gamma(p)$  consists of two consecutive parts,  $\gamma_0(p)$  and  $\Gamma_1(p)$ . The path  $\gamma_0(p)$  is contained in  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$  and the path  $\Gamma_1(p)$  is contained in  $\text{Zer}(Q, \mathbb{R}^k)_{p_1}$ . The part  $\gamma_0(p)$  consists of a sequence of sub-paths,  $\gamma_{0,0}, \dots, \gamma_{0,m}$ . Each  $\gamma_{0,i}$  is a semi-algebraic path parametrized by one of the co-ordinates  $X_1, \dots, X_k$ , over some interval  $[a_{0,i}, b_{0,i}]$  with  $\gamma_{0,0}(a_{0,0}) = p$ . The semi-algebraic maps,  $\gamma_{0,0}, \dots, \gamma_{0,m}$  and the end-points of their intervals of definition  $a_{0,0}, b_{0,0}, \dots, a_{0,m}, b_{0,m}$  are all independent of  $p$  (up to the discrete choice of the path  $\gamma(p)$  in  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \{p\})$ ), except  $b_{0,m}$  which depends on  $p_1$ .

Moreover,  $\Gamma_1(p)$  can again be decomposed into two parts,  $\gamma_1(p)$  and  $\Gamma_2(p)$  with  $\Gamma_2(p)$  contained in  $\text{Zer}(Q, \mathbb{R}^k)_{p_2}$  and so on.

If  $q = (q_1, \dots, q_k) \in \text{Zer}(Q, \mathbb{R}^k)$  is another point such that  $p_1 \neq q_1$ , then since  $\text{Zer}(Q, \mathbb{R}^k)_{p_1}$  and  $\text{Zer}(Q, \mathbb{R}^k)_{q_1}$  are disjoint, it is clear that

$$\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \{p\}) \cap \text{RM}(\text{Zer}(Q, \mathbb{R}^k), \{q\}) = \text{RM}(\text{Zer}(Q, \mathbb{R}^k)).$$

Now consider a connecting path  $\gamma(q)$  extracted from  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \{q\})$ . The images of  $\Gamma_1(p)$  and  $\Gamma_1(q)$  are disjoint. If the image of  $\gamma_0(q)$  (which is contained in  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$ ) follows the same sequence of curve segments as  $\gamma_0(p)$  starting at  $p$  (that is, it consists of the same curves segments  $\gamma_{0,0}, \dots, \gamma_{0,m}$  as in  $\gamma_0(p)$ ), then it is clear that the images of the paths  $\gamma(p)$  and  $\gamma(q)$  has the property that they are identical up to a point and they are disjoint after it. We call this the **divergence property**.  $\square$

Next we show how to handle the case when the input algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  is not bounded.

*Algorithm 15.5. [Algebraic Roadmap]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$  together with a finite set  $\mathcal{N}$  of real univariate representations with coefficients in  $D$ .
- **Output:** a roadmap  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \mathcal{N})$  which contains  $\mathcal{N}$ .
- **Complexity:**  $d^{O(k^2)}$ , where  $d$  is a bound on the degree of  $Q$  and  $O(d)^k$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ , the degrees of the univariate representations in  $\mathcal{N}$ , and the number of these univariate representations.
- **Procedure:**
  - Introduce new variables  $X_{k+1}$  and  $\varepsilon$  and replace  $Q$  by the polynomial

$$Q_\varepsilon = Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_{k+1}^2) - 1)^2.$$

Replace  $\mathcal{N} \subset \mathbb{R}^k$  by the set of real univariate representations specifying the elements of  $\text{Zer}(\varepsilon^2 (X_1^2 + \dots + X_{k+1}^2) - 1, \mathbb{R}\langle\varepsilon\rangle^{k+1})$  above the points associated to  $\mathcal{N}$  using Algorithm 12.11 (Univariate Representation).

- Run Algorithm 15.3 (Bounded Algebraic Roadmap) without a triangular Thom encoding (i.e. with  $i = 0$ ),  $Q_\varepsilon$  and  $\mathcal{N}$  as input with structure  $D[\varepsilon]$ . The algorithm outputs a roadmap of  $\text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}), \mathcal{N})$  composed of points and curves whose description involves  $\varepsilon$ .
- Denote by  $\mathcal{L}$  the set of polynomials in  $D[\varepsilon]$  whose signs have been determined in the preceding computation and take  $a = \min_{P \in \mathcal{L}} c'(P)$  (Definition 10.5). Replace  $\varepsilon$  by  $a$  in the polynomial  $Q_\varepsilon$  to get a polynomial  $Q_a$ . Replace  $\varepsilon$  by  $a$  in the output roadmap to obtain a roadmap  $\text{RM}(\text{Zer}(Q_a, \mathbb{R}^{k+1}), \mathcal{N})$ . When projected to  $\mathbb{R}^k$ , this roadmap gives a roadmap for  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \mathcal{N}) \cap B(0, 1/a)$ .
- In order to extend the roadmap outside the ball  $B(0, 1/a)$  collect all the points  $(y_1, \dots, y_k, y_{k+1}) \in \mathbb{R}\langle\varepsilon\rangle^{k+1}$  in the roadmap  $\text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}), \mathcal{N})$  which satisfies  $\varepsilon(y_1^2 + \dots + y_k^2) = 1$ . Each such point is described by a real univariate representation involving  $\varepsilon$ . Add to the roadmap the curve segment obtained by first forgetting the last coordinate and then treating  $\varepsilon$  as a parameter which varies over  $(0, a, ]$  to get a roadmap  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \mathcal{N})$ .

**Proof of correctness:** The choice of  $a$  guarantees that the roadmap for  $Q_\varepsilon$  just computed specializes to a roadmap for  $Q_a$  when  $\varepsilon$  is replaced by  $a$ . The correctness follows from the correctness of Algorithm 15.3 (Bounded Algebraic Roadmap). □

**Complexity analysis:** According to the complexity analysis of Algorithm 15.3 (Bounded Algebraic Roadmap), the number of arithmetic operations in the ring  $D[\varepsilon]$  is  $d^{O(k^2)}$ . Moreover, the degrees of the polynomials in  $\varepsilon$  generated by the algorithm do not exceed  $d^{O(k^2)}$ , using the complexity analysis of Algorithm 12.10 (Parametrized Special Multiplication Table). The complexity is thus  $d^{O(k^2)}$  in the ring  $D$ , taking into account the complexity analyses of Algorithm 8.4 (Addition of multivariate polynomials), Algorithm 8.5 (Multiplication of Multivariate Polynomials), and Algorithm 8.6 (Exact Division of Multivariate Polynomials).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

*Algorithm 15.6. [Algebraic Connecting]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$ ,
  - a real univariate representation  $u, \tau$  with coefficients in  $D$ , with associated point  $p$  contained in  $\text{Zer}(Q, R^k)$ .
- **Output:** a path  $\gamma(p, \text{Zer}(Q, R^k)) \subset \text{Zer}(Q, R^k)$  connecting  $p$  to a distinguished point of  $\text{RM}(\text{Zer}(Q, R^k))$ .
- **Complexity:**  $d^{O(k^2)}$ , where  $d$  is a bound on the degree of  $Q$  and  $O(d)^k$  is a bound on the degrees of  $u$ .
- **Procedure:** Call Algorithm 15.5 (Algebraic Roadmap) with input  $Q$  and  $(u, \tau)$  and extract  $\gamma$  from  $\text{RM}(\text{Zer}(Q, R^k), \{u, \tau\})$ .

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 15.5 (Algebraic Roadmap). □

**Complexity analysis:** The total number of arithmetic operations in  $D$  is bounded by  $d^{O(k^2)}$ , using the complexity analysis of Algorithm 15.5 (Algebraic Roadmap).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

We can now summarize our results on the complexity of the computation of the roadmap for an algebraic set.

**Theorem 15.11.** *Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$  be a polynomial whose total degree is at most  $d$ .*

- a) *There is an algorithm whose output is exactly one point in every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ . The complexity in the ring generated by the coefficients of  $Q$  is bounded by  $d^{O(k^2)}$ . In particular, this algorithm counts the number of semi-algebraically connected components of  $\text{Zer}(Q, \mathbb{R}^k)$  in time  $d^{O(k^2)}$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .*
- b) *Let  $p$  and  $q$  in  $\text{Zer}(Q, \mathbb{R}^k)$  be two points which are represented by real  $k$ -univariate real representation  $u, \sigma v, \tau$  of degree  $O(d)^k$ . There is an algorithm deciding whether  $p$  and  $q$  belong to the same connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ . The complexity in the ring generated by the coefficients of  $Q$  and the coefficients of the polynomials in  $u$  and  $v$  is bounded by  $d^{O(k^2)}$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .*

**Proof:** For a), proceed as follows: first compute  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$ , then describe its connected components using the adjacencies between curve segments and points, and finally take one point in each of these connected components.

For b), use Algorithm 15.6 (Algebraic Connecting) for  $p$  and  $q$ . The points  $p$  and  $q$  are connected to points  $p'$  and  $q'$  of the roadmap. Use the first item to decide whether they belong to the same connected component or not.  $\square$

### 15.3 Computing Connected Components of Algebraic Sets

This section is devoted to the proof of the following result.

**Theorem 15.12.** *If  $\text{Zer}(Q, \mathbb{R}^k)$  is an algebraic set defined as the zero set of a polynomial  $Q \in \mathbb{D}[X_1, \dots, X_k]$  of degree  $\leq d$ , then there is an algorithm that outputs quantifier free formulas whose realizations are the semi-algebraically connected components of  $\text{Zer}(Q, \mathbb{R}^k)$ . The complexity of the algorithm in the ring generated by the coefficients of  $Q$  is bounded by  $d^{O(k^3)}$  and the degrees of the polynomials that appear in the output are bounded by  $O(d)^{k^2}$ . Moreover, if  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^3)}$ .*

The proof is based on a parametrized version of the roadmap algorithm: we are going to find sign conditions on the parameters for which the description of the roadmap does not change.

For this purpose, we need parametrized versions of Algorithm 12.21 (Triangular Comparison of Roots) and Algorithm 12.22 (Triangular Intermediate Points). These algorithms will be based on Algorithm 14.6 (Parametrized Sign Determination).

Let  $\mathcal{A} \subset \mathcal{B}$ ,  $\rho$  and  $\bar{\rho}$  two sign conditions on  $\mathcal{A}$  and  $\mathcal{B}$ . The sign condition  $\bar{\rho}$  **refines**  $\rho$  if  $\bar{\rho}(P) = \rho(P)$  for every  $P \in \mathcal{A}$ .

**Notation 15.13.** We denote by  $\text{SIGN}(\rho, \mathcal{B})$  the list of realizable sign conditions on  $\mathcal{B}$  refining  $\rho$ .  $\square$

*Algorithm 15.7.* [**Parametrized Comparison of Roots**]

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .
- **Input:** a parametrized Thom encoding  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$ , of level  $k - 1$ , with coefficients in  $D$ , two non-zero polynomials  $P$  and  $Q \in D[Y, X_1, \dots, X_k]$ .
- **Output:**
  - a finite set  $\mathcal{B} \subset D[Y]$  containing  $\mathcal{A}$ ,
  - for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B})$ , a list of sign conditions on  $\text{Der}(\mathcal{T} \cup \{P\} \cup \{Q\})$  refining  $\sigma$  specifying for every  $y \in \text{Reali}(\rho)$  the ordered list of the triangular Thom encodings of the roots of  $P$  and  $Q$  above the point specified by  $\sigma$ .
- **Complexity:**  $d^{O(k\ell)}$ , where  $\ell$  is the number of parameters and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ , and the degree of  $P$  and  $Q$ .
- **Procedure:** Apply Algorithm 14.6 (Parametrized Sign Determination) to  $\mathcal{T}$ ,  $P$  and

$$\text{Der}(\mathcal{T}) \cup \text{Der}(P) \cup \text{Der}(Q),$$

then to  $\mathcal{T}$ ,  $Q$  and

$$\text{Der}(\mathcal{T}) \cup \text{Der}(P) \cup \text{Der}(Q)$$

**Proof of correctness:** Immediate.  $\square$

**Complexity analysis:** The complexity is  $d^{O(k\ell)}$ , using the complexity of Algorithm 14.6 (Parametrized Sign Determination). The number of elements in  $\mathcal{B}$  is  $d^{O(k\ell)}$ , and the degrees of the elements of  $\mathcal{A}$  are bounded by  $d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k\ell)}$ .  $\square$

*Algorithm 15.8.* [**Parametrized Intermediate Points**]

- **Structure:** an ordered integral domain  $D$  contained in a real closed field  $R$ .



- **Input:** a parametrized Thom encoding  $\mathcal{A}, \rho, \mathcal{T}, \sigma$  of level  $k - 1$ , with coefficients in  $D$ , two non-zero polynomials  $P$  and  $Q$  in  $D[Y, X_1, \dots, X_k]$  of degree bounded by  $p$ .
- **Output:**
  - a finite set  $\mathcal{B} \subset D[Y]$  containing  $\mathcal{A}$
  - for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B})$ , a list of sign conditions on  $\text{Der}(\mathcal{T} \cup \{(PQ)'\})$  specifying for every  $y \in \text{Reali}(\bar{\rho})$  the triangular Thom encodings of a set of points intersecting all the intervals between two consecutive roots of  $P$  and  $Q$ .
- **Complexity:**  $d^{O(k\ell)}$ , where  $\ell$  is the number of parameters and  $d$  is a bound on the degrees of the polynomials in  $\mathcal{T}$ , and the degree of  $P$  and  $Q$ .
- **Procedure:** Apply Algorithm 14.7 (Parametrized Thom Encoding) with input  $\mathcal{T}, P, \mathcal{T}, Q$  and  $\mathcal{T}, P'Q$ . Sort them using Algorithm 15.7 (Parametrized Comparison of Roots).

**Proof of correctness:** Immediate. □

**Complexity analysis:** The complexity is  $d^{O(k\ell)}$ , using the complexity of Algorithm 14.6 (Parametrized Sign Determination). The number of elements in  $\mathcal{A}$  is  $d^{O(k\ell)}$ , and the degrees of the elements of  $\mathcal{B}$  are bounded by  $d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k\ell)}$ . □

A **parametrized real univariate triangular representation** of level  $i - 1$  with parameters  $Y = (Y_1, \dots, Y_\ell)$   $\mathcal{T}, \sigma, u$  above  $\mathcal{A}, \rho$  is

- a parametrized triangular Thom encoding  $\mathcal{T}, \sigma$  of level  $i$ ,
- a parametrized representation  $u = (\mathcal{T}_i, g_0, g_i, \dots, g_k) \subset D[Y, X_{\leq i}, T]$  such that for every  $y \in \text{Reali}(\rho)$  there is a root  $(z(y), t(y))$  of  $\mathcal{T}$  with triangular Thom encoding  $\sigma$ .

A parametrized real univariate triangular representation  $\mathcal{T}, \sigma, u$  is **above** the parametrized triangular Thom encoding  $\mathcal{A}, \rho, \mathcal{T}', \sigma'$  if  $\mathcal{T}, \sigma$  is above  $\mathcal{A}, \rho$  and if  $\mathcal{T}' = \mathcal{T}_1, \dots, \mathcal{T}_{i-1}$ , and  $\sigma' = \sigma_1, \dots, \sigma_{i-1}$ .

*Algorithm 15.9. [Parametrized Projection]*

- **Structure:** a domain  $D$  contained in a field  $K$ .
- **Input:** a parametrized real univariate representation  $\mathcal{T}, u, \sigma$  above a parametrized triangular Thom encoding  $\mathcal{A}, \rho$ , with coefficients in  $D$ . For every  $y \in \text{Reali}(\rho)$ , we denote by  $z(y)$  the root of  $\mathcal{T}(y)$  specified by  $\tau$  and by  $x(y)$  the point associated to  $u(y, z(y))$ .
- **Output:**
  - a finite set  $\mathcal{B} \subset D[Y]$  containing  $\mathcal{A}$ ,

- for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B})$  a Thom encoding  $(\text{proj}_i(u), \text{proj}_i(\tau))$  specifying, for every  $y \in \text{Reali}(\bar{\rho})$ , the projection of the point associated to  $x(y)$  on the  $X_i$  axis.

- **Complexity:**  $d^{O(ki\ell)}$ , where  $\ell$  is the number of parameters,  $d$  is a bound on the degrees of on the degree of the univariate representation and of the polynomials in  $\mathcal{T}$ .

- **Procedure:**

- Compute the resultant  $\text{proj}_i(u)$  of  $f(Y, X_{<i}, T)$ , and

$$X_i g_0(Y, X_{<i}, T) - g_i(Y, X_{<i}, T)$$

with respect to  $T$ , using Algorithm 8.21 (Signed subresultant).

- Use Algorithm 14.6 (Parametrized Sign Determination) with  $\mathcal{T}$ ,  $f$  and the derivatives of

$$g_0(Y, X_{<i}, T)^d \text{proj}_i(u) \left( \frac{g_i(Y, X_{<i}, T)}{g_0(Y, X_{<i}, T)} \right)$$

with respect to  $T$ , where  $d$  is the smallest even number not less than the degree of  $\text{proj}_i(u)$  with respect to  $X_i$ . This gives a list of polynomials  $\mathcal{B} \subset D[Y]$  and for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B})$  the Thom encoding  $\text{proj}_i(\tau)$  of the  $i$ -th coordinate of  $x(y)$ .

**Proof of correctness:** Immediate. □

**Complexity analysis:** The complexity is  $d^{O(ki\ell)}$ , using the complexity of Algorithm 14.6 (Parametrized Sign Determination).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(ki\ell)}$ . □

We now define parametrized curve segments.

Let  $\mathcal{V}_1, \tau_1, \mathcal{V}_2, \tau_2$  be two parametrized triangular Thom encoding above  $\mathcal{A}, \rho, \mathcal{T}, \sigma$ . For every  $y \in \text{Reali}(\rho)$ , we denote by  $z(y) \in \mathbb{R}^{i-1}$  the point specified by  $\mathcal{T}(y), \sigma$  and by  $(z(y), a(y)), (z(y), b(y))$  the points specified by  $\mathcal{V}_1(y), \tau_1$  and  $\mathcal{V}_2(y), \tau_2$ . A **parametrized curve segment representation**  $u, \tau$  above  $\mathcal{V}_1, \tau_1, \mathcal{V}_2, \tau_2$  is given by

- a parametrized univariate representation with parameters  $(Y, X_{\leq i})$ ,

$$u = (f(Y, X_{\leq i}, T), g_0(Y, X_{\leq i}, T), g_{i+1}(Y, X_{\leq i}, T), \dots, g_k(Y, X_{\leq i}, T)),$$

- a sign condition  $\tau$  on  $\text{Der}(f)$  such that for every  $y \in \text{Reali}(\rho)$  and for every  $v \in (a(y), b(y))$  there exists a real root  $t(v)$  of  $f(z(y), v, T)$  with Thom encoding  $\sigma, \rho, \tau$  and  $g_0(z(y), v, t(v)) \neq 0$ .

Our aim is first to describe a parametrized version of Algorithm 15.2 (Curve Segments).

**Algorithm 15.10. [Parametrized Curve Segments]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a parametrized Thom encoding  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with parameters  $Y = (Y_1, \dots, Y_\ell)$  of level  $i - 1$ , with coefficients in  $D$ . For every  $y \in \text{Reali}(\rho)$ ,  $(y, z(y))$  denotes the point specified by  $\sigma$ .
  - a polynomial  $Q \in D[Y, X_1, \dots, X_k]$ , for which  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$
  - a finite set  $\mathcal{N}$  of parametrized real univariate triangular representation above  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with, for every  $y \in \text{Reali}(\rho)$ , associated points contained in  $\text{Zer}(Q, R^k)$ .

• **Output:**

- a finite set  $\mathcal{B} \subset D[Y]$  containing  $\mathcal{A}$ ,
- for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B})$ ,
  - an ordered list of parametrized Thom encodings

$$\mathcal{V}_{\bar{\rho}, 1}, \tau_{\bar{\rho}, 1}, \dots, \mathcal{V}_{\bar{\rho}, \ell(\bar{\rho})}, \tau_{\bar{\rho}, \ell(\bar{\rho})}$$

above  $\mathcal{B}$ ,  $\bar{\rho}$ ,  $\mathcal{T}$ ,  $\sigma$

- for every  $i = 1, \dots, \ell(\bar{\rho})$ ,
  - a finite set  $\mathcal{N}_{\bar{\rho}, i}$  of parametrized real univariate triangular representations above

$$\mathcal{B}, \bar{\rho}, \mathcal{V}_{\bar{\rho}, j}, \tau_{\bar{\rho}, j}$$

- a finite set  $\mathcal{C}_{\bar{\rho}, j}$  of parametrized curve segments above

$$\mathcal{B}, \bar{\rho}, \mathcal{V}_{\bar{\rho}, j}, \tau_{\bar{\rho}, j}, \mathcal{V}_{\bar{\rho}, j+1}, \tau_{\bar{\rho}, j+1}$$

- a list of pairs of elements of  $\mathcal{C}_{\bar{\rho}, j}$  and  $\mathcal{N}_{\bar{\rho}, j}$  (resp.  $\mathcal{C}_{\bar{\rho}, j+1}$  and  $\mathcal{N}_{\bar{\rho}, j}$ ) describing the adjacency relation.

For every  $y \in \text{Reali}(\bar{\rho})$ , this defines a set of curves and points contained in  $\text{Zer}(Q, R^k)_{y, z(y)}$ . The specifications of these points and curves is fixed for every point  $y \in \text{Reali}(\bar{\rho})$ . These points and curves satisfy the properties of the output of Algorithm 15.2 (Curve Segments).

- **Complexity:**  $d^{O(ki\ell)}$ , where  $\ell$  is the number of parameters,  $d$  is a bound on the degree of  $Q$ ,  $O(d)^k$  is a bound on the degrees of on the degree of the parametrized univariate representations in  $\mathcal{N}$  and of the polynomials in  $\mathcal{T}$ .

• **Procedure:**

- Step 1: Perform Algorithm 12.10 (Parametrized Multiplication Table) with input  $\overline{\text{Cr}}(Q^2, \zeta)$ , using Notation 12.46, and parameter  $Y$ ,  $X_{\leq i}$ . Perform Algorithm 12.15 (Parametrized Limit of Bounded Points), and output a set  $\mathcal{U}$  of parametrized univariate representations.

Using Notation 13.8, consider for every  $u = (f, g_0, g_{i+1}, \dots, g_k) \in \mathcal{U}$  the finite set  $\mathcal{F}_u$  containing  $Q_u$  and all the derivatives of  $f$  with respect to  $T$ , and compute  $\mathcal{D}_u = \text{RElim}_T(f, \mathcal{F}_u) \subset D[Y, X_{\leq i}]$  using Algorithm 11.19 (Restricted Elimination).

- Define  $\mathcal{D} = \bigcup_{u \in \mathcal{U}} \mathcal{D}_u$ .

- Step 2: Use Algorithm 15.9 (Parametrized Projection) with input  $\mathcal{N}$  and output a finite set  $\mathcal{B}_2 \subset \mathbb{D}[Y]$  containing  $\mathcal{A}$ , such that for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B}_2)$  and every  $u \in \mathcal{N}$  the Thom encoding  $\text{proj}_i(u)$ ,  $\text{proj}_i(\tau)$  specifying the projection of the associated point on the  $X_i$  axis is fixed for every  $y \in \text{Reali}(\bar{\rho})$ . Add to  $\mathcal{D}$  the polynomials  $\text{proj}_i(u)$ .
- Step 3: Apply Algorithms 14.7 (Parametrized Thom Encoding), 15.7 (Parametrized Comparison of Roots) to the set  $\mathcal{D}$ . Denote by  $\mathcal{B}_3 \subset \mathbb{D}[Y]$  the family of polynomials output, and for every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B}_3)$ , denote by

$$A_{\bar{\rho},1}\alpha_{\bar{\rho},1}, \dots, A_{\bar{\rho},\ell(\bar{\rho})}\alpha_{\bar{\rho},\ell(\bar{\rho})}$$

the list of Thom encodings output. For every  $y \in \text{Reali}(\bar{\rho})$ , these are the Thom encodings of the corresponding distinguished values

$$v_1(y, z(y)) < \dots < v_\ell(y, z(y)).$$

Define  $\mathcal{V}_i, \tau_i = \mathcal{T}, A_i$  and  $\tau_i = \sigma, \alpha_i$ .

- Step 4: For every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B}_3)$ , every  $j = 1, \dots, \ell(\bar{\rho})$  and every  $u = (f, g_0, g_i, \dots, g_k), \tau \in \mathcal{N}$ , use Algorithm 14.7 (Parametrized Triangular Thom Encoding) and output  $\mathcal{B}_4(\bar{\rho}, j, u)$ , containing  $\mathcal{B}_3$ . Append pairs  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{N}_{\rho_1,j}$  for every  $\rho_1 \in \text{SIGN}(\bar{\rho}, \mathcal{B}_4(\bar{\rho}, j, u, \tau))$  such that for every  $y \in \text{Reali}(\rho_1)$   $\text{proj}_i(\tau)$  is the Thom encoding of a point of  $\text{Zer}(Q, \mathbb{R}^k)_z(y)$  with projection having Thom encoding  $\alpha_j$ . Define  $\mathcal{B}_4(\bar{\rho}) = \cup \mathcal{B}_4(\bar{\rho}, j, u, \tau)$ .
- Step 5: For every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B}_3)$ , every  $j = 1, \dots, \ell(\bar{\rho})$  and every  $u = (f, g_0, g_i, \dots, g_k) \in \mathcal{U}$ , use Algorithm 15.8 (Parametrized Intermediate Points) and Algorithm 14.7 (Parametrized Triangular Thom Encoding) and output  $\mathcal{B}_5(\bar{\rho}, j, u)$ , containing  $\mathcal{B}_3$ . Append pairs  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{N}_{\rho_1,j}$  for every  $\rho_1 \in \text{SIGN}(\bar{\rho}, \mathcal{B}_5(\bar{\rho}, j, u))$  such that for every  $y \in \text{Reali}(\rho_1)$   $\text{proj}_i(\tau)$  is the Thom encoding of a point of  $\text{Zer}(Q, \mathbb{R}^k)_z(y)$  with projection having Thom encoding  $\alpha_j$ . Define  $\mathcal{B}_5(\bar{\rho}) = \cup \mathcal{B}_5(\bar{\rho}, j, u)$ .
- Step 6: For every  $\bar{\rho} \in \text{SIGN}(\rho, \mathcal{B}_3)$ , every  $j = 1, \dots, \ell(\bar{\rho}) - 1$  and every  $u = (f, g_0, g_i, \dots, g_k) \in \mathcal{U}$ , use Algorithm 15.8 (Parametrized Intermediate Points) and Algorithm 14.7 (Parametrized Triangular Thom Encoding) and output a family  $\mathcal{B}_6(\bar{\rho}, j, u)$  containing  $\mathcal{B}_3$  such that for every sign condition  $\rho_1$  on  $\mathcal{B}_6$  and every  $y \in \text{Reali}(\rho_1)$  the Thom encodings  $\tau$  of the roots of  $f(y, z(y), v, T)$  over  $(v_i(y), v_{i+1}(y))$  are fixed and the corresponding associated curves are contained in  $\text{Zer}(Q, \mathbb{R}^k)_z(y)$ . Append all pairs  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{C}_{\rho_3,i}$ . Define  $\mathcal{B}_6(\bar{\rho}) = \cup \mathcal{B}_6(\bar{\rho}, j, u)$ .
- Step 7: Consider  $\rho_1 \in \text{SIGN}(\bar{\rho}, \mathcal{B}_4 \cup \mathcal{B}_5 \cup \mathcal{B}_6)$ . For every  $j = 1, \dots, \ell(\bar{\rho}_1)$  and every parametrized real univariate triangular representation of  $\mathcal{N}_{\rho_1,j}$  specified by

$$v' = (p, q_0, q_2, \dots, q_k), \tau', \{p, q_0, q_2, \dots, q_k\} \subset \mathbb{D}[Y, X_{\leq i}][T]$$

and every parametrized curve segment representation of  $\mathcal{C}_{\rho_1, j}$  specified by

$$v = (f, g_0, g_2, \dots, g_k), \tau, \{f, g_0, g_2, \dots, g_k\} \subset D[Y, X_{\leq i}[T],$$

compute a family  $\mathcal{B}_7(\rho_1, v', \tau', v, \tau)$  of polynomials containing  $\mathcal{B}_4 \cup \mathcal{B}_5 \cup \mathcal{B}_6$  such that for every  $\rho_2 \in \text{SIGN}(\rho_1, \mathcal{B}_7(\rho_1, v', \tau', v, \tau))$  and every  $y \in \text{Reali}(\rho_2)$  the algorithm deciding whether the corresponding point  $t(y)$  is adjacent to the corresponding curve segment gives the same answer: compute the first  $\nu$  such that  $(\partial^\nu g_0 / \partial X_i^\nu)(v_j, t)$  is not zero and decide whether for every  $\ell = i + 1, \dots, k$

$$\frac{\partial^\nu g_\ell}{\partial X_i^\nu}(v_j, t)q_0(t) - \frac{\partial^\nu g_0}{\partial X_i^\nu}(v_j, t)q_\ell(t)$$

is zero, using Algorithm 14.6 (Parametrized Sign Determination).

Repeat the same process for every element of  $\mathcal{N}_{\rho_1, i+1}$  and every curve segment of  $\mathcal{C}_{\rho_1, i}$ .

- Finally output  $\mathcal{B} = \cup \mathcal{B}_7(\rho_1, v', \tau', v, \tau)$ .

**Proof of correctness:** It follows from Proposition 12.42 and the correctness of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination), Algorithm 15.9 (Parametrized Projection), Algorithm 15.8 (Parametrized Intermediate Points), Algorithm 14.7 (Parametrized Thom Encoding), Algorithm 15.7 (Parametrized Comparison of Roots) and Algorithm 14.6 (Parametrized Sign Determination).  $\square$

### Complexity analysis:

- Step 1: This step requires  $d^{O((\ell+i)(k-i))}$  arithmetic operations in  $D$ , using the complexity analyses of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination). There are  $d^{O(k-i)}$  parametrized univariate representations computed in this step and each polynomial in these representations has degree  $O(d)^{k-i}$ .
- Step 2: This step requires  $d^{O((\ell+i)k)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 15.9 (Parametrized Projection).
- Step 3: This step requires  $d^{O(\ell ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 14.7 (Parametrized Thom Encoding).
- Step 4: This step requires  $d^{O(\ell ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 14.6 (Parametrized Sign Determination).
- Step 5: This step requires  $d^{O(\ell ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 14.7 (Parametrized Thom Encoding).
- Step 6: This step requires  $d^{O(\ell ik)}$  arithmetic operations in  $D$ , using the complexity analyses of Algorithm 15.8 (Parametrized Intermediate Points) and Algorithm 14.7 (Parametrized Thom Encoding).

- Step 7: This step requires  $d^{O(\ell ik)}$  arithmetic operations, using the complexity analysis of Algorithm 14.6 (Parametrized Sign Determination).

Thus, the complexity is  $d^{O(\ell ik)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(\ell ik)}$ .  $\square$

*Algorithm 15.11.* **[Parametrized Bounded Algebraic Roadmap]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a parametrized Thom encoding  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with parameters  $Y = (Y_1, \dots, Y_\ell)$  and variables  $X_{\leq i} = (X_1, \dots, X_i)$ , with coefficients in  $D$ . For every  $y \in \text{Reali}(\rho)$ ,  $(y, z(y))$  denotes the point specified by  $\sigma$ ,
  - a polynomial  $Q \in D[Y, X_1, \dots, X_k]$ , for which  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a finite set  $\mathcal{N}$  of parametrized real univariate triangular representations above  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with coefficients in  $D$ , with, for every  $y \in \text{Reali}(\rho)$ , associated points contained in  $\text{Zer}(Q, \mathbb{R}^k)$ .
- **Output:**
  - a subset  $\mathcal{C}$  of  $D[Y]$  containing  $\mathcal{A}$ ,
  - for every realizable sign condition  $\tau$  on  $\mathcal{C}$  refining  $\rho$ , a subset  $\text{RM}(\tau)$  such that, for every  $y \in \text{Reali}(\tau)$ ,  $\text{RM}(\tau)_y$  is a roadmap for  $\text{Zer}(Q, \mathbb{R}^k)_y$  that contains  $\mathcal{N}_y$ .
- **Complexity:**  $d^{O(\ell k^2)}$ , where  $\ell$  is the number of parameters,  $O(d)^k$  is a bound on the degrees of on the degree of the univariate representation and of the polynomials in  $\mathcal{T}$ .
- **Procedure:**
  - Call Algorithm 15.10 (Parametrized Curve Segments), output  $\mathcal{B}$  and, for every realizable sign condition  $\bar{\rho}$  on  $\mathcal{B}$  refining  $\rho$ ,  $\ell(\bar{\rho})$ . Output also, for every  $j = 1, \dots, \ell(\bar{\rho})$ ,  $A_{\bar{\rho}, j}$ ,  $\alpha_{\bar{\rho}, j}$ ,  $\mathcal{N}_{\bar{\rho}, j}$  and  $\mathcal{C}_{\bar{\rho}, j}$ .
  - For every realizable sign condition  $\bar{\rho}$  on  $\mathcal{B}$  and for every  $i$  from 1 to  $\ell(\bar{\rho})$ , call Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap) recursively, with input  $\mathcal{B}$ ,  $\bar{\rho}$ ,  $\mathcal{T}$ ,  $A_{\bar{\rho}, j}$ ,  $\sigma$ ,  $\alpha_{\bar{\rho}, j}$ ,  $Q$  and  $\mathcal{N}_{\bar{\rho}, j}$ .

**Proof of correctness:** The correctness of the algorithm follows from Proposition 15.7 and the correctness of Algorithm 15.2 (Curve Segments).  $\square$

**Complexity analysis:** In the recursive calls to Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap), the number of triangular systems considered is at most  $d^{O(k^2)}$  and the triangular systems involved have polynomials of degree  $O(d)^k$ .

Thus, the total number of arithmetic operations in  $D$  is bounded by  $d^{O(\ell k^2)}$  using the complexity analysis of Algorithm 15.10 (Parametrized Curve Segments).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(\ell k^2)}$ .  $\square$

We now want to obtain a parametrized connecting algorithm. We show how to obtain a covering of a given  $\mathcal{P}$ -closed semi-algebraic set contained in  $\text{Zer}(Q, \mathbb{R}^k)$  by a family of semi-algebraically contractible subsets. The construction is based on a parametrized version of the connecting algorithm: we compute a family of polynomials such that for each realizable sign condition  $\sigma$  on this family, the description of the connecting paths of different points in the realization,  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$ , are uniform.

We first define parametrized paths. A parametrized path is a semi-algebraic set which is a union of semi-algebraic paths having the divergence property (see Remark 15.10).

More precisely,

**Definition 15.14.** A **parametrized path**  $\gamma$  is a continuous semi-algebraic mapping from  $V \subset \mathbb{R}^{k+1} \rightarrow \mathbb{R}^k$ , such that, denoting by  $U = \pi_{1\dots k}(V) \subset \mathbb{R}^k$ , there exists a semi-algebraic continuous function  $\ell: U \rightarrow [0, +\infty)$ , and there exists a point  $a$  in  $\mathbb{R}^k$ , such that

- $V = \{(x, t) \mid x \in U, 0 \leq t \leq \ell(x)\}$ ,
- $\forall x \in U, \gamma(x, 0) = a$ ,
- $\forall x \in U, \gamma(x, \ell(x)) = x$ ,
- $\forall x \in U, \forall y \in U, \forall s, 0 \leq s \leq \ell(x), \forall t, 0 \leq t \leq \ell(y)$   
 $(\gamma(x, s) = \gamma(y, t) \Rightarrow s = t)$ ,
- $\forall x \in U, \forall y \in U, \forall s \in [0, \min(\ell(x), \ell(y))]$   
 $(\gamma(x, s) = \gamma(y, s) \Rightarrow \forall t \leq s \gamma(x, t) = \gamma(y, t))$ .  $\square$

Given a parametrized path,  $\gamma: V \rightarrow \mathbb{R}^k$ , we will refer to  $U = \pi_{1\dots k}(V)$  as its *base*. Also, any semi-algebraic subset  $U' \subset U$  of the base of such a parametrized path, defines in a natural way the restriction of  $\gamma$  to the base  $U'$ , which is another parametrized path, obtained by restricting  $\gamma$  to the set  $V' \subset V$ , defined by  $V' = \{(x, t) \mid x \in U', 0 \leq t \leq \ell(x)\}$ .

**Proposition 15.15.** *Let  $\gamma: V \rightarrow \mathbb{R}^k$  be a parametrized path such that  $U = \pi_{1\dots k}(V)$  is closed and bounded. Then, the image of  $\gamma$  is semi-algebraically contractible.*

**Proof:** Let  $W = \text{Im}(\gamma)$  and  $M = \sup_{x \in U} \ell(x)$ . We prove that the semi-algebraic mapping  $\phi: W \times [0, M] \rightarrow W$  sending

$$\begin{aligned} &(\gamma(x, t), s) \text{ to } \gamma(x, s) \text{ if } t \geq s, \\ &(\gamma(x, t), s) \text{ to } \gamma(x, t) \text{ if } t < s \end{aligned}$$

is continuous. Note that the map  $\phi$  is well-defined, since

$$\gamma(x, t) = \gamma(x', t') \Rightarrow t = t',$$

by condition (4). Since  $\phi$  satisfies

$$\begin{aligned} \phi(\gamma(x, t), 0) &= a, \\ \phi(\gamma(x, t), M) &= \gamma(x, t) \end{aligned}$$

this gives a semi-algebraic continuous contraction from  $W$  to  $\{a\}$ .

Let  $w \in W, s \in [0, M]$ . Let  $\varepsilon > 0$  be an infinitesimal, and let

$$(w', s') \in \text{Ext}(W \times [0, M], \mathbb{R}\langle\varepsilon\rangle)$$

be such that  $\lim_\varepsilon (w', s') = (w, s)$ . In order to prove the continuity of  $\phi$  at  $w$  it suffices to prove that

$$\lim_\varepsilon \text{Ext}(\phi, \mathbb{R}\langle\varepsilon\rangle)(w', s') = \phi(w, s).$$

Let  $w = \gamma(x, t)$  for some  $x \in U, t \in [0, \ell(x)]$ , and similarly let  $w' = \gamma(x', t')$  for some  $x' \in \text{Ext}(U, \mathbb{R}\langle\varepsilon\rangle)$  and  $t' \in [0, \text{Ext}(\ell, \mathbb{R}\langle\varepsilon\rangle)(x')]$ . Note that  $\lim_\varepsilon (x') \in U$  since  $U$  is closed and bounded and  $\lim_\varepsilon t' \in [0, \ell(\lim_\varepsilon x')]$ .

Now,

$$\begin{aligned} \gamma(x, t) &= w \\ &= \lim_\varepsilon (w') \\ &= \lim_\varepsilon \text{Ext}(\gamma, \mathbb{R}\langle\varepsilon\rangle)(x', t') \\ &= \gamma(\lim_\varepsilon x', \lim_\varepsilon t'). \end{aligned}$$

Condition (4) now implies that  $\lim_\varepsilon t' = t$ .

Without loss of generality let  $t' \geq t$ . The other case is symmetric. We have the following two sub-cases.

- Case  $s' > t'$ : Since  $s, t \in \mathbb{R}$  and  $\lim_\varepsilon s' = s$  and  $\lim_\varepsilon t' = t$ , we must have that  $s \geq t$ . In this case  $\text{Ext}(\phi, \mathbb{R}\langle\varepsilon\rangle)(w', s') = \text{Ext}(\gamma, \mathbb{R}\langle\varepsilon\rangle)(x', t')$ . Then,

$$\begin{aligned} \lim_\varepsilon \text{Ext}(\phi, \mathbb{R}\langle\varepsilon\rangle)(w', s') &= \lim_\varepsilon \text{Ext}(\gamma, \mathbb{R}\langle\varepsilon\rangle)(x', t') \\ &= \lim_\varepsilon w' \\ &= w \\ &= \phi(w, s). \end{aligned}$$

- Case  $s' \leq t'$ : Again, since  $s, t \in \mathbb{R}$  and  $\lim_\varepsilon s' = s$  and  $\lim_\varepsilon t' = t$ , we must have that  $s \leq t$ .

In this case we have,

$$\begin{aligned} \lim_\varepsilon \phi(w', s') &= \lim_\varepsilon \text{Ext}(\gamma, \mathbb{R}\langle\varepsilon\rangle)(x', s') \\ &= \gamma(\lim_\varepsilon x', \lim_\varepsilon s') \\ &= \gamma(\lim_\varepsilon x', s). \end{aligned}$$



Now,

$$\begin{aligned}
 \gamma(\lim_{\varepsilon} x', t) &= \gamma(\lim_{\varepsilon} x', \lim_{\varepsilon} t') \\
 &= \lim_{\varepsilon} \text{Ext}(\gamma, \mathbf{R}(\varepsilon))(x', t') \\
 &= \lim_{\varepsilon} w' \\
 &= w \\
 &= \gamma(x, t).
 \end{aligned}$$

Thus, by condition (5) we have that  $\gamma(\lim_{\varepsilon} x', s'') = \gamma(x, s'')$  for all  $s'' \leq t$ . Since,  $s \leq t$ , this implies,

$$\begin{aligned}
 \lim_{\varepsilon} \text{Ext}(\phi, \mathbf{R}(\varepsilon))(w', s') &= \lim_{\varepsilon} \text{Ext}(\gamma, \mathbf{R}(\varepsilon))(w', s') \\
 &= \gamma(\lim_{\varepsilon} x', \lim_{\varepsilon} s') \\
 &= \gamma(x, s) \\
 &= \phi(w, s).
 \end{aligned}$$

This proves the continuity of  $\phi$ , using Proposition 3.5.  $\square$

*Algorithm 15.12.* **[Parametrized Bounded Algebraic Connecting]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a parametrized Thom encoding  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with parameters  $Y = (Y_1, \dots, Y_{\ell})$  and variables  $X_{\leq i} = (X_1, \dots, X_i)$ , with coefficients in  $D$ . For every  $y \in \text{Reali}(\rho)$ ,  $(y, z(y))$  denotes the point specified by  $\sigma$ ,
  - a polynomial  $Q \in D[Y, X_1, \dots, X_k]$ , for which  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$
  - a parametrized real univariate triangular representation above  $\mathcal{A}$ ,  $\rho$ ,  $\mathcal{T}$ ,  $\sigma$  with coefficients in  $D$ , with, for every  $y \in \text{Reali}(\rho)$ , associated point  $p(y)$  contained in  $\text{Zer}(Q, R^k)$ .
- **Output:**
  - a subset  $\mathcal{C}$  of  $D[Y]$  containing  $\mathcal{A}$ ,
  - for every realizable sign condition  $\tau$  on  $\mathcal{C}$  refining  $\rho$ , a parametrized path  $\gamma(\tau)$  such that, for every  $y \in \text{Reali}(\tau)$ ,  $\gamma(\tau)(y)$  is a path connecting  $p(y)$  to a distinguished point of  $\text{RM}(\text{Zer}(Q, R^k))$ .
- **Complexity:**  $d^{O(\ell k^2)}$ , where  $\ell$  is the number of parameters,  $O(d)^k$  is a bound on the degrees of on the degree of the univariate representation and of the polynomials in  $\mathcal{T}$ .
- **Procedure:** Call Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap) and extract  $\gamma$  from  $\text{RM}(\tau)$ .

**Proof of correctness:** The correctness of the algorithm follows from the correctness of Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap). It is easy to see that  $\gamma$  is a parametrized path (see Definition 15.14), using the divergence property of the paths  $\gamma(y, \cdot)$  (see Remark 15.10).  $\square$

**Complexity analysis:** The total number of arithmetic operations in  $D$  is bounded by  $d^{O(\ell k^2)}$ , using the complexity analysis of Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(\ell k^2)}$ .  $\square$

*Algorithm 15.13.* **[Connected Components of an Algebraic Set]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$ .
- **Output:** a subset  $\mathcal{A}$  of  $D[X_1, \dots, X_k]$  and for every semi-algebraically connected component  $S$  of  $\text{Zer}(Q, \mathbb{R}^k)$  a finite subset  $\Sigma \subset \text{SIGN}(\mathcal{A})$  such that  $S = \bigcup_{\sigma \in \Sigma} \text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$ .
- **Complexity:**  $d^{O(k^3)}$ , where  $d$  is a bound on the degree of the polynomial  $Q$ .
- **Procedure:**
  - Take  $Q_\varepsilon = Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2$ .
  - Call Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap) without parametrized triangular Thom encoding,  $Q_\varepsilon$ , and

$$\mathcal{N} = \{(T - 1, 1, Y_1, \dots, Y_k)\}.$$

The output contains a family of polynomials  $\mathcal{A}^* \subset D[\varepsilon][X]$  such that the realization of a non-empty sign condition  $\rho$  in  $\mathcal{A}^*$  is contained in a semi-algebraically connected component of  $\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1})$ .

- Find a set  $\mathcal{S}$  of sample points for every realizable sign condition on  $\mathcal{A}^*$  using Algorithm 13.1 (Sampling). Compute  $\text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}))$  using Algorithm 15.3 (Bounded Algebraic Roadmap) and for every semi-algebraically connected component  $S'$  of  $\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1})$ , fix a point  $y(S')$  of  $S' \cap \text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}))$ . For every  $x \in \mathcal{S}$  compute a roadmap  $\text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}), x)$  of  $\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1})$  containing  $x$  using Algorithm 15.3 (Bounded Algebraic Roadmap) and decide from  $\text{RM}(\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1}), x)$  whether  $x$  belongs to  $S'$ .
- Output the description of  $S'$ , i.e. the disjunction  $\Phi(S')$  of realizable sign conditions on  $\mathcal{A}^*$  with a sample point belonging to  $S'$ , for every semi-algebraically connected component  $S'$  of  $\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1})$ .
- For every connected component  $S$  of  $\text{Zer}(Q, \mathbb{R}^k)$  there exists a connected component  $S'$  of  $\text{Zer}(Q_\varepsilon, \mathbb{R}\langle\varepsilon\rangle^{k+1})$ , such that  $\pi(S') \cap \mathbb{R}^k = S$ , where  $\pi: \mathbb{R}\langle\varepsilon\rangle^{k+1} \rightarrow \mathbb{R}\langle\varepsilon\rangle^k$  is the projection map forgetting the last coordinate.

Consider the formula  $\Phi(S')$  describing  $S'$  and, eliminating a quantifier, the formula  $\Psi$  describing  $\pi(S')$ . Then  $\text{Remo}_\varepsilon(\Psi(Y))$  (Notation 14.6) defines  $S$ .

**Proof of correctness:** All points satisfying the same sign condition on  $\mathcal{A}$  can be connected by a semi-algebraic path in  $\text{Zer}(Q, \mathbb{R}^k)$  to some fixed curve segment of  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$  and hence must belong to the same connected component of  $\text{Zer}(Q, \mathbb{R}^k)$ . Which realizable sign conditions on  $\mathcal{A}$  belong to the same semi-algebraically connected component of  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$  follows from Step 2 and 3. We also use Proposition 14.7.  $\square$

**Complexity analysis:** The total number of arithmetic operations in  $\mathbb{D}$  is bounded by  $d^{O(k^3)}$ , using the complexity analysis of Algorithm 15.11 (Parametrized Bounded Algebraic Roadmap). The degrees of the polynomials in  $\mathcal{A}$  are bounded by  $d^{O(k^2)}$ .

If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^3)}$ .  $\square$

So we have proved Theorem 15.12.

## 15.4 Bibliographical Notes

The problem of deciding connectivity properties of algebraic sets considered here is a base case for deciding connectivity properties of semi-algebraic sets, studied in Chapter 16.

The notion of a roadmap for a semi-algebraic set was introduced by Canny in [36].

We discuss in more details the various contributions to the roadmap problem and the computation of connected components at the end of Chapter 16.

It is interesting to remark that the complexity of computing the number of connected components of an algebraic set given in this chapter is significantly worse than that of the algorithm for computing the Euler-Poincaré characteristic of an algebraic set given in Chapter 12. Thus, currently we are able to compute the Euler-Poincaré characteristic of real algebraic sets (which is the alternative sum of the Betti numbers) more efficiently than any of the individual Betti numbers.

---

## Computing Roadmaps and Connected Components of Semi-algebraic Sets

We compute roadmaps and connected components of semi-algebraic sets. The algorithms described in this chapter have complexity much better than the ones provided by cylindrical decomposition in Chapter 11 for the problem of deciding connectivity properties of semi-algebraic sets (single exponential in the number of variables rather than doubly exponential).

In Section 16.2, we study uniform roadmaps, which provide roadmaps in the realization of every weak sign condition obtained by the relaxation of a realizable sign condition of a finite set of polynomials  $\mathcal{P}$ . A key algorithm is the Connecting Algorithm which links a point to the uniform roadmap inside the same weak sign condition. The correctness of the uniform roadmap algorithm relies on properties of special values studied in Section 16.1.

In Section 16.3, using a parametrized version of the Connecting Algorithm we show how to compute descriptions of the semi-algebraically connected components of the realizations of sign conditions by quantifier free formulas.

In Section 16.4, we show how to compute descriptions of the semi-algebraically connected components of semi-algebraic sets by quantifier free formulas. In Section 16.5 we construct roadmaps for general semi-algebraic sets. Finally in Section 16.6 we give a single exponential complexity algorithm for computing the first Betti number of a semi-algebraic set.

### 16.1 Special Values

We want to prove a result similar to Proposition 15.4 for basic semi-algebraic sets. Unfortunately, the notion of pseudo-critical values is not strong enough to ensure this property, and this is why we define the technical notion of special value.

Let  $\text{Zer}(Q, \mathbb{R}^k)$  be bounded. Suppose that

- $Q \in \mathbb{D}[X_1, \dots, X_k]$ , of degree at most  $d$  is such that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
- $d_1 \geq d_2 \geq \dots \geq d_k$ ,
- $\deg(Q) \leq d_1$ ,  $\text{tDeg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ .

Let  $\bar{d}_i = 2d_i + 2$ ,  $i = 1, \dots, k$ , and  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ . Consider

$$\text{Def}(Q^2, \zeta) = \zeta G_k(\bar{d}, c) + (1 - \zeta) Q^2,$$

using Notation 12.46.

An  $X_1$ -**special value** of  $\text{Zer}(Q, \mathbb{R}^k)$  is a  $c \in \mathbb{R}$  for which there exists  $y \in \text{Zer}(\text{Def}(Q^2, \bar{d}, c, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  with  $\lim_{\zeta} (\pi(y)) = c$ ,  $g(y)$  infinitesimal and  $y$  a local minimum of  $g$  on  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$ , where

$$g(X) = \frac{\sum_{i=2}^k \left( \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_i} \right)^2}{\sum_{i=1}^k \left( \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_i} \right)^2} \tag{16.1}$$

Note that any  $X_1$ -pseudo-critical value of  $\text{Zer}(Q, \mathbb{R}^k)$  is an  $X_1$ -special value of  $\text{Zer}(Q, \mathbb{R}^k)$ .

Let  $S$  be a basic closed semi-algebraic set defined as

$$S = \{x \in \mathbb{R}^k \mid Q(x) = 0 \wedge \bigwedge_{P \in \mathcal{P}} P(x) \geq 0\}.$$

An  $X_1$ -**special value on  $S$**  is an  $X_1$ -special value on  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$  where  $\mathcal{P}'$  is contained in  $\{Q\} \cup \mathcal{P}$ .

Using special values, a result similar to Proposition 15.4 holds for basic closed semi-algebraic sets.

**Proposition 16.1.** *Let  $\text{Zer}(Q, \mathbb{R}^k)$  be bounded, and let  $S$  be a basic closed semi-algebraic set defined as*

$$S = \{x \in \mathbb{R}^k \mid Q(x) = 0 \wedge \bigwedge_{P \in \mathcal{P}} P(x) \geq 0\}.$$

*If  $C$  is a semi-algebraic connected component of  $S_{[a,b]}$  and  $[a,b] \setminus \{v\}$  contains no  $X_1$ -special value of  $S$ ,  $v \in (a,b)$ , then  $C_v$  is semi-algebraically connected.*

We first prove the following

**Proposition 16.2.** *If  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$  and  $x$  is a point of  $\text{Zer}(Q, \mathbb{R}^k)_v$  at which  $\text{Zer}(Q, \mathbb{R}^k) \cap B(x, \varepsilon)_{<v}$  is empty for a positive  $\varepsilon$ , then  $v$  is a  $X_1$ -special value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof:** We first prove that the statement of the proposition can be translated into a formula of the language of ordered fields with coefficients in  $\mathbb{R}$ . More precisely, we prove that the statement

$$\exists y \in \text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k) \quad \lim_{\zeta} (\pi(y)) = v \wedge \lim_{\zeta} (g(y)) = 0 \tag{16.2}$$

is equivalent to the formula

$$\forall \varepsilon > 0 \quad \exists \delta > 0 \quad \forall t \quad 0 < t < \delta \quad \exists y \quad \Phi(v),$$

with

$$\Phi(v) := \text{Def}(Q^2, t)(y) = 0 \wedge g_t(y)^2 + (\pi(y) - v)^2 < \varepsilon,$$

where we write  $g_t$  for the rational fraction obtained after replacing  $\zeta$  by  $t$  in the definition of  $g$ . Note that, for  $t$  small enough,  $g_t$  is well defined and the values of  $g_t$  are bounded by 1. Thus the limit  $g_0$  of  $g_t$ , as  $t$  tends to 0, is well defined.

First observe that Equation (16.2) is equivalent to the fact that there exists a germ  $\varphi$  of semi-algebraic function represented by a semi-algebraic continuous function  $h$  defined on  $[0, a]$  such that  $\pi(h(0)) = v$  and  $g_0(h(0)) = 0$  by Corollary 3.11, Proposition 3.18 and Lemma 3.21, since  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle\zeta\rangle^k)$  is bounded by Proposition 12.38. Equation (16.2) follows from the continuity of  $h$ , taking  $y = h(t)$ .

In the other direction, assume

$$\forall \varepsilon > 0 \quad \exists \delta > 0 \quad \forall t \quad 0 < t < \delta \quad \exists y \quad \Phi(v).$$

By Theorem 3.19 (Curve selection lemma), there exists  $\varepsilon_0 > 0$  and a semi-algebraic continuous function  $d$  from  $[0, \varepsilon_0]$  to  $\mathbb{R}^k$  such that  $d(0) = 0$  and for all  $0 < \varepsilon < \varepsilon_0$ ,  $d(\varepsilon) > 0$  and

$$\forall t \quad 0 < t < d(\varepsilon) \quad \exists y \quad \text{Def}(Q^2, t)(y) = 0 \wedge g_t(y)^2 + (\pi(y) - v)^2 < \varepsilon. \tag{16.3}$$

Since  $d(\varepsilon_0) \in \mathbb{R}$  is positive, we can find, using Proposition 3.4,  $\varepsilon$  infinitesimal in  $\mathbb{R}\langle\zeta\rangle$  such that  $\text{Ext}(d, \mathbb{R}\langle\zeta\rangle)(\varepsilon) = 2\zeta$ . Choosing  $t = \zeta$ , there exists  $y \in \text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle\zeta\rangle^k)$  such that  $g(y)$  and  $\pi(y) - v$  are infinitesimal. This proves Equation (16.2).

Considering all polynomials  $Q$  of fixed degree, the statement of the proposition can now be expressed by a sentence of the language of ordered fields with coefficients in  $\mathbb{Z}$ . By Theorem 2.80 (Tarski-Seidenberg principle), it thus suffices to prove the proposition over the reals, which is what we now proceed to do. The proposition for  $\mathbb{R} = \mathbb{R}$  is an immediate consequence of the following two lemmas.

Let  $g$  be defined in Equation (16.1).

**Lemma 16.3.** *Suppose that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$  and that  $x$  is a point of  $\text{Zer}(Q, \mathbb{R}^k)_v$  at which  $\text{Zer}(Q, \mathbb{R}^k) \cap B(x, \varepsilon)_{<v}$  is empty for some positive  $\varepsilon$ , then there is a point  $y \in \text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle\zeta\rangle^k) \cap B(x, \varepsilon)$  for which  $\lim_\zeta (\pi(y)) = v$  and  $\lim_\zeta (g(y)) = 0$ .*

**Lemma 16.4.** *If  $y$  is a point of  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle\zeta\rangle^k) \cap B(x, \varepsilon)$  at which  $\lim_\zeta (\pi(y)) = v$  and  $\lim_\zeta (g(y)) = 0$  then  $v$  is a  $X_1$ -special value of  $\pi$  on  $\text{Zer}(Q, \mathbb{R}^k)$ .*

**Proof of Lemma 16.3 :** If there is a critical value of  $\pi$  on

$$\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle\zeta\rangle^k)$$

infinitesimally close to  $v$ , we are done. Otherwise, suppose that there is no critical value of  $\pi$  on  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  in an interval  $(v - b, v + b) \subset \mathbb{R}\langle \zeta \rangle$  with  $b \in \mathbb{R}$ . We can suppose without loss of generality that  $b > \varepsilon$ .

We argue by contradiction and suppose that for every  $y$  at which

$$\text{Def}(Q^2, \zeta)(y) = 0 \wedge \lim_{\zeta} (\pi(y)) = c,$$

the value  $g(y)$  is not infinitesimal.

Since  $\text{Zer}(Q, \mathbb{R}^k) \cap B(x, \varepsilon)_{<c} = \emptyset$ , we know that for any

$$y \in \text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k) \cap B(x, \varepsilon)_{\leq v},$$

$\lim_{\zeta} (\pi(y)) = v$  and thus  $g(y)$  is not infinitesimal. Let  $a \in \mathbb{R}$  be a positive number smaller than any value of  $g$  on

$$\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k) \cap B(x, \varepsilon)_{\leq v}.$$

Let

$$U' = \{t \in \mathbb{R} \mid g_t < a \text{ on } \text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k) \cap B(x, \varepsilon)_{\leq v}\}.$$

Let  $U''$  be the set of  $t \in \mathbb{R}$  such that there is no critical value of  $\pi$  on  $\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k)$  in  $(v - b, v + b)$  and  $U = U' \cap U''$ . The set  $U$  is semi-algebraic and its extension to  $\mathbb{R}\langle \zeta \rangle$  contains  $\zeta$ . Thus, it contains an interval  $(0, t_0)$  by Proposition 3.17.

For every  $t \in (0, t_0)$ , let  $y_t$  be a point in

$$\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k) \cap B(x, \varepsilon)_{\leq c}$$

whose last  $k - 1$  coordinates coincide with the last  $k - 1$  coordinates of  $x$ . Consider the curve  $\gamma_t$  on  $\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k)$  through  $y_t$  which at each of its points is tangent to the gradient of  $\pi$  on  $\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k)$ . The gradient of  $\pi$  on  $\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k)$  at a point of  $\text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k)$  is proportional to

$$G = \left( \sum_{i=2}^k \left( \frac{\partial \text{Def}(Q^2, t)}{\partial X_i} \right)^2, \dots, -\frac{\partial \text{Def}(Q^2, t)}{\partial X_1} \frac{\partial \text{Def}(Q^2, t)}{\partial X_k} \right)$$

(see page 240). For every point of  $\gamma_t$ , the vector  $G$  thus belongs to the half-cone  $\mathcal{C}$  of center  $x$ , based on the  $k - 1$ -sphere of radius  $\sqrt{\frac{1-a}{a}}$  and center  $(x_1 - 1, x_2, \dots, x_k)$  in the hyperplane  $X_1 = x_1 - 1$ . It follows that the curve  $\gamma_t$  is completely contained in  $\mathcal{C}$ . Since there is no critical value of  $\pi$  on  $\text{Zer}(Q_t, \mathbb{R}^k)$  in  $(v - b, v + b)$ , the curve  $\gamma_t$  is defined over  $(v - b, v + b)$  and thus meets  $S(x, \varepsilon) \cap \mathcal{C}$ .

Since  $\mathcal{C} \cap S(x, \varepsilon) \cap \text{Zer}(\text{Def}(Q^2, t), \mathbb{R}^k) \neq \emptyset$  is true for every  $t \in (0, t_0)$  it follows from Proposition 3.17 that

$$\mathcal{C} \cap S(x, \varepsilon) \cap \text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k) \neq \emptyset.$$

Thus, taking  $\lim_{\zeta}$  of the point so obtained,  $B(x, \varepsilon)_{<v} \cap \text{Zer}(Q, \mathbb{R}^k) \neq \emptyset$ , which is a contradiction. □

**Proof of Lemma 16.4 :** If  $g$  is zero anywhere that the first coordinate is infinitesimally close to  $v$ , then  $v$  is a  $X_1$ -pseudo-critical value and we are done. Alternatively, we may assume that  $g$  is non-zero in any slab of infinitesimal width containing  $X_1 = v$ . Let  $y$  be given by our hypothesis, i.e.  $\lim_{\zeta} (\pi(y)) = v$ ,  $\lim_{\zeta} (g(y)) = 0$ . We let  $C$  be the bounded semi-algebraically connected component of  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  containing  $y$ . Define  $w$  by  $\pi(y) = w$ . Then  $g$  attains its minimum on  $C_w$  at some point  $z \in C_w$ . Let  $t$  be this minimum. It is clear that  $t$  is infinitesimal.

Consider the set  $A = \{w \mid \min_{C_w} (g) \leq t\}$ . This set  $A$  is closed, bounded, semi-algebraic, and thus a union of closed intervals  $[a_1, b_1] \cup \dots \cup [a_h, b_h]$  with  $a_i \leq b_i < a_{i+1}$ . Let  $[a_i, b_i] = [a, b]$  be the interval containing  $w$ .

If  $a$  and  $b$  are both infinitesimally close to  $w$  take  $u$  and  $u'$  so that  $b_{i-1} < u < a = a_i \leq b = b_i < u' < a_{i+1}$  with  $u$  and  $u'$  infinitesimally close to  $w$ . The minimum of  $g$  on  $C_{[u, u']}$  occurs in the interior of the slab since it is smaller at  $C_w$  than its minimum both on  $C_u$  and  $C_{u'}$ . It follows that  $c$  is a  $X_1$ -special value on  $\text{Zer}(Q, \mathbb{R}^k)$ .

Assume on the contrary that  $[a, b]$  is such that  $a$  or  $b$  is not infinitesimally close to  $w$ . We are going to prove that this leads to a contradiction.

According to Theorem 5.46 (Semi-algebraic triviality), there exists a family  $\phi_j$  of semi-algebraic curves parametrized by open segments  $(\alpha_j, \beta_j)$  covering  $(a, b)$  (with the exception of a finite number of points) such that  $g(\phi_j(x))$  is smaller than  $t$ . If  $T_j(x) = (T_{j,1}(x), \dots, T_{j,k}(x))$  is the tangent vector to  $\phi_j$  at  $(x, \phi_j(x))$ , we have

$$\begin{aligned}
 -T_{j,1} \frac{\partial \text{Def}(Q^2, t)}{\partial X_1} &= T_{j,2} \frac{\partial \text{Def}(Q^2, t)}{\partial X_2} + \dots + T_{j,k} \frac{\partial \text{Def}(Q^2, t)}{\partial X_k} \\
 T_{j,1}^2 &\leq \frac{t}{1-t} \|(T_{j,2}, \dots, T_{j,k})\|^2.
 \end{aligned}$$

Thus, at every point on each of these curves,  $\left| \frac{T_{j,1}(x)}{T_{j,i}(x)} \right| < \sqrt{\frac{kt}{1-t}} = t'$  for some  $2 \leq i \leq k$ . Hence, we can suppose – subdividing further if needed and producing more curves – that on each of these curves,  $\left| \frac{T_{j,1}(x)}{T_{j,i}(x)} \right| < t'$  for some  $2 \leq i \leq k$ .

Let  $N$  be the number of the curves so obtained. We prove now that the interval  $(w, w + 2/cNt')$  contains  $w$  such that  $\min_{C_w} (g) > t$ . Suppose on the contrary that at every value  $u \in (w, w + 2/cNt')$ ,  $\min_{C_u} (g) \leq t$ . Then there is an interval of length at least  $2/ct'$  over which the curve  $\phi_j(x)$  is differentiable and  $\left| \frac{T_{j,1}(x)}{T_{j,i}(x)} \right|$  is less than  $t'$ . It follows from the mean value theorem that the projection of this curve to the  $X_i$  axis is bigger than  $2/c$ , which contradicts the fact that  $C \subset B(0, 1/c)$ . Similarly, the interval  $(w - 2/cNt', w)$  contains  $u'$  such that  $\min_{C_{u'}} (g) > t$ .



Note that both  $u$  and  $u'$  are infinitesimally close to  $v$ . This contradicts the fact that  $a$  or  $b$  is not infinitesimally close to  $w$  and ends the argument.  $\square$   $\square$

The proof of Proposition 16.1 will use the following lemma.

Let  $C$  be a semi-algebraically connected component of  $S_{[a,v]}$  and let  $B_1, \dots, B_h$  be the semi-algebraically connected components of  $C_{[a,v]}$ .

**Lemma 16.5.** *If  $\bar{B}_1 \cap \bar{B}_2 \neq \emptyset$ , then  $v$  is a  $X_1$ -pseudo-critical value on  $S$ .*

**Proof:** Suppose that  $\bar{B}_1 \cap \dots \cap \bar{B}_I \neq \emptyset$  and that  $1, \dots, I$  is a maximal family with this property. Let  $x$  be a point of this intersection. Clearly,  $x$  belongs to the boundary of  $S$  and the set  $\mathcal{P}' \subset \mathcal{P}$  of polynomials in  $\mathcal{P}$  that vanish at  $x$  is not empty. According to Theorem 5.46 (Semi-algebraic triviality) there is  $w \in [a, v]$  such that  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)_{[w,v]}$  is semi-algebraically homeomorphic to  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)_w \times [w, v]$  and  $C_{[w,v]}$  is semi-algebraically homeomorphic to  $C_w \times [w, v]$ . Note that  $C_{[w,v]}$  is not semi-algebraically connected. Let  $D$  be the connected component of  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)_{[w,v]}$  containing  $x$ .

We consider two cases according to whether or not  $D_w$  is empty:

If  $D_w$  is empty, then  $v$  is an  $X_1$ -pseudo-critical value on  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$  by Proposition 16.2 and we have already noted that pseudo-critical values are special values.

If  $D_w$  is not empty, then some semi-algebraically connected component of  $C_{[a,v]}$  intersects  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$  in any neighborhood of  $x$ . Suppose, without loss of generality that it is  $B_1$ . Consider a maximal subset of  $\mathcal{P}$ , say  $\mathcal{P}''$ , such that  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)$  intersects  $B_2$  in any neighborhood of  $x$ . The set  $\mathcal{P}''$  is non-empty and contained in  $\mathcal{P}'$ . According to Theorem 5.46 (Semi-algebraic triviality) there is a  $w' \geq w$  such that  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)_{[w',v]}$  is semi-algebraically homeomorphic to  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)_{w'} \times [w', v]$ . Let  $Z$  be the connected component of  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)_{[w',v]}$  containing  $x$ . By the maximality of  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)$ , there is a connected component  $Z_1$  of  $Z_{[w',v]}$  contained in  $B_{2[w',v]}$ . Since  $\text{Zer}(\mathcal{P}', \mathbb{R}^k) \subset \text{Zer}(\mathcal{P}'', \mathbb{R}^k)$  and  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)_{[w',v]}$  meets  $B_1$ ,  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)_{[w',v]}$  is not semi-algebraically connected. We conclude by Proposition 15.5 that  $v$  is a  $X_1$ -pseudo-critical value on  $\text{Zer}(\mathcal{P}'', \mathbb{R}^k)$ .  $\square$

**Proof of Proposition 16.1:** Suppose that  $C_v$  is empty. We take  $d \in [a, b]$  such that  $C_d$  is non-empty and suppose that  $v < d$  (the case  $v > d$  can be treated similarly). We obtain a contradiction by proving that there is a  $X_1$ -special value on  $S$  in  $(v, d]$ . Since the set  $\{w \in (v, d] \mid C_w \neq \emptyset\}$  is a closed semi-algebraic subset of  $[v, d]$ , it contains a smallest such value, say  $u$ . Choose an  $x \in C_u$ . Since  $x$  belongs to the boundary of  $S$ , the set  $\mathcal{P}'$  of polynomials in  $\mathcal{P}$  vanishing at  $x$  is non-empty. It is clear that  $\text{Zer}(\mathcal{P}', \mathbb{R}^k) \cap B(x, \varepsilon)_{<u} = \emptyset$  for  $\varepsilon$  small enough. Hence, by Proposition 16.2,  $u$  is an  $X_1$ -special value on  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$ .

Suppose now that  $C_v$  is not semi-algebraically connected. Take  $d \in [a, b]$  such that a semi-algebraically connected component of  $C_{[v,d]}$  contains more than one connected component of  $C_v$  and suppose that  $v < d$  (the case  $v > d$  can be treated similarly). We obtain a contradiction by proving that there is a  $X_1$ -special value on  $S$  in  $(v, d]$ . Since the set of  $w \in (v, d]$  for which  $C_{[v,w]}$  contains more than one connected component of  $C_v$  is a closed semi-algebraic subset of  $[v, b]$  by Theorem 5.46 (Semi-algebraic triviality), it contains a smallest such value, say  $u$ .

Consider a connected component  $B$  of  $C_{[v,u]}$  containing more than one connected component of  $C_u$ . Let  $B_1, \dots, B_h$  be the connected components of  $C_{[v,u]}$  contained in  $B$ , and let  $B_0$  be the set of  $x \in B_u$  such that  $B(x, \varepsilon)_{<u} \cap C = \emptyset$  for  $\varepsilon$  small enough. Clearly,  $B = B_0 \cup \overline{B_1} \cup \dots \cup \overline{B_h}$ .

We now prove that  $u$  is an  $X_1$ -special value on  $S$  whether or not  $B_0 = \emptyset$ .

If  $B_0$  is non-empty, choose an  $x \in B_0$  and let  $\mathcal{P}'$  be the set of polynomials in  $\mathcal{P}$  vanishing at  $x$ . Then  $B(x, \varepsilon)_{<u} \cap \text{Zer}(\mathcal{P}', \mathbb{R}^k)$  is empty and it follows from Proposition 16.2 that  $u$  is an  $X_1$ -special value on  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$ .

Alternatively, if  $B_0$  is empty we may assume, without loss of generality, that  $\overline{B_1} \cap \overline{B_2} \neq \emptyset$ . Thus by Lemma 16.5,  $u$  is an  $X_1$ -pseudo-critical value, hence a  $X_1$ -special value on  $S$ . □

We are going now to indicate how to compute special values. Consider the algebraic set  $Z$  defined by the  $k + 1$  polynomial equations in the  $k + 1$  variables  $(X_1, \dots, X_k, \lambda)$

$$\begin{aligned} \text{Def}(Q^2, \zeta) &= 0, \\ \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_1} &= \lambda \frac{\partial g}{\partial X_1}, \\ &\vdots \\ \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_k} &= \lambda \frac{\partial g}{\partial X_k}. \end{aligned}$$

The local minima of  $g$  on  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  are contained in the projection of  $Z$  to the first  $k$  coordinates.

**Proposition 16.6.** *If  $C'$  is a semi-algebraically connected component of  $Z$  on which  $g$  has an infinitesimal local minimum on  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  then  $\lim_{\zeta} (C')$  is a single point.*

**Proof:** Let  $x$  be a point of  $C'$  where  $g$  has an infinitesimal local minimum on  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$ , and let  $u = g(x)$ . Note that  $g$  is constant on  $C'$ . The projection of  $C'$  to the  $X_1$ -axis,  $\pi(C')$ , is contained in  $A = \{w \mid \min_{C_w} (g) \leq u\}$  where  $C$  is the semi-algebraically connected component of  $\text{Zer}(\text{Def}(Q^2, \zeta), \mathbb{R}\langle \zeta \rangle^k)$  containing  $x$ . Since  $\pi(C')$  is semi-algebraically connected, following the proof of Lemma 16.4 we see that  $\pi(C')$  is contained in an infinitesimal segment. □

*Algorithm 16.1.* [Special Values]

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a triangular system  $\mathcal{T}$  specifying  $z \in R^{i-1}$ , with coefficients in  $D$ ,
  - a polynomial  $Q \in D[X_1, \dots, X_k]$ , such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$ .
- **Output:** a set of values containing the  $X_i$ -special values of  $\text{Zer}(Q, R^k)_z$ .
- **Complexity:**  $d^{O(k)}$  where  $d$  is the degree of  $Q$ .
- **Procedure:**
  - Let  $d_1 \geq d_2 \dots \geq d_k$ ,  $\deg(Q) \leq d_1$ ,  $\text{tDeg}_{X_i}(Q) \leq d_i$ , for  $i = 2, \dots, k$ ,  $\bar{d}_i = 2d_i + 2$ ,  $i = 1, \dots, k$ ,  $\bar{d} = (\bar{d}_1, \dots, \bar{d}_k)$ , and

$$\text{Def}(Q^2, \zeta) = \zeta G_{k-i}(\bar{d}, c) + (1 - \zeta) Q^2,$$

using Notation 12.46. Denote by  $Z$  the algebraic set defined by the  $k + 1$  polynomial equations in the  $k + 1$  variables  $(X_1, \dots, X_k, \lambda)$ ,

$$\begin{aligned} \mathcal{T}_j(X_1, \dots, X_j) &= 0, \mathcal{T}_j \in \mathcal{T}, j = 1, \dots, i - 1, \\ \text{Def}(Q^2, \zeta) &= 0, \\ \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_i} &= \lambda \frac{\partial g}{\partial X_i}, \\ &\vdots \\ \frac{\partial \text{Def}(Q^2, \zeta)}{\partial X_k} &= \lambda \frac{\partial g}{\partial X_k}. \end{aligned}$$

- Use Algorithm 12.16 (Bounded Algebraic Sampling) to find a set of points which meets every connected component of  $Z$ .
- Compute  $\lim_\zeta$  of these points using Algorithm 12.14 (Limit of Real Bounded Points).
- Describe their  $k$  first coordinates using Algorithm 15.1 (Projection). Keep only the points whose  $i - 1$  first coordinates coincide with  $z$ .

**Proof of correctness:** By Proposition 16.6, we know that the  $X_i$ -special values of  $\text{Zer}(Q, R^k)_z$  are the among the values computed by Algorithm 16.1 (Special Values). □

**Complexity analysis:** Using the complexity analysis of Algorithm 12.16 (Bounded Algebraic Sampling), Algorithm 12.14 (Limit of Real Bounded Points), and Algorithm 15.1 (Projection), we conclude that the complexity is  $d^{O(k)}$ . At most  $O(d)^k$  univariate polynomials of degrees at most  $O(d)^k$  whose real roots contain the  $X_1$ -special values of  $\text{Zer}(Q, R^k)$  are computed.

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ . □

## 16.2 Uniform Roadmaps

We consider

- a polynomial  $Q \in \mathbb{R}[X_1, \dots, X_k]$  for which  $\text{Zer}(Q, \mathbb{R}^k)$  is bounded,
- a set of at most  $s$  polynomials  $\mathcal{P}$  such that  $\mathcal{P}$  is in strong  $\ell$ -general position with respect to  $Q$  (Definition b).

We first indicate how to connect any point  $x \in \text{Zer}(Q, \mathbb{R}^k)$  to some roadmap of the zero set of the union of  $Q$  and a subset of  $\mathcal{P}$ .

Denote by  $\sigma(x)$  the sign condition on  $\mathcal{P}$  at  $x$ . Let

$$\text{Reali}(\bar{\sigma}(x), \text{Zer}(Q, \mathbb{R}^k)) = \{x \in \text{Zer}(Q, \mathbb{R}^k) \mid \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) \in \bar{\sigma}(x)(P)\},$$

where  $\bar{\sigma}$  is the relaxation of  $\sigma$  (Definition 5.32). We say that  $\bar{\sigma}(x)$  is the weak sign condition defined by  $x$  on  $\mathcal{P}$ . We denote by  $\mathcal{P}(x)$  the union of  $\{Q\}$  and the set of polynomials in  $\mathcal{P}$  vanishing at  $x$ .

### Algorithm 16.2. [Bounded Connecting]

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a finite set of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$  in strong  $\ell$ -general position with respect to  $Q$ ,
  - a point  $p \in \text{Zer}(Q, \mathbb{R}^k)$  described by a real univariate representation  $u, \tau$  with coefficients in  $D$ .
- **Output:** a subset  $\mathcal{P}' \subset \mathcal{P}$  and a semi-algebraic path  $\Gamma$  which connects  $p$  to  $\text{RM}(\text{Zer}(\mathcal{P}' \cup \{Q\}, \mathbb{R}^k))$  inside  $\text{Reali}(\bar{\sigma}(x), \text{Zer}(Q, \mathbb{R}^k))$ .
- **Complexity:**  $\ell s d^{O(k^2)}$ , where  $s$  is a bound on the number of polynomials in  $\mathcal{P}$ ,  $d$  is a bound on the degree of  $Q$  and the polynomials in  $\mathcal{P}$  and  $O(d)^k$  is a bound on the degree of of the univariate representation  $u$ .
- **Procedure:**
  - Initialize  $\Gamma = \emptyset, q := p, u', \tau' := u, \tau$ .
  - $(\star)$  Construct a path  $\gamma$  connecting  $q$  to  $\text{RM}(\text{Zer}(\mathcal{P}(q), \mathbb{R}^k), \{u', \tau'\})$ , using Algorithm 15.4 (Bounded Algebraic Connecting).
  - If a polynomial of  $\mathcal{P} \setminus \mathcal{P}(q)$  vanishes somewhere on  $\gamma$ , let  $t \in (0, 1)$  such that no polynomial in  $\mathcal{P} \setminus \mathcal{P}(q)$  vanishes on  $\gamma((0, t))$  and there are polynomials in  $\mathcal{P} \setminus \mathcal{P}(q)$  vanishing on  $\gamma(t)$ . Add  $\gamma|_{(0, t]}$  to the end of  $\Gamma$ . Call the algorithm recursively returning to  $(\star)$  with input  $q := \gamma(t)$ , taking as  $u', \tau'$  the real univariate representation describing  $\gamma(t)$ .
  - If  $\gamma$  is such that no polynomial of  $\mathcal{P} \setminus \mathcal{P}(q)$  vanishes on  $\gamma$ , add  $\gamma$  to the end of  $\Gamma$ .

**Proof of correctness :** Follows clear from the correctness of Algorithm 15.4 (Bounded Algebraic Connecting).  $\square$

**Complexity analysis:** Since  $\mathcal{P}$  is in strong  $\ell$ -general position with respect to  $Q$ , the algorithm terminates after  $\ell' \leq \ell$  iterations. The degrees of the univariate representations representing the  $\ell'$  successive values of  $p$  are bounded by  $d^{O(k^2)}$ . Thus the complexity of the Bounded Connecting Algorithm is clearly  $\ell s d^{O(k^2)}$ . The number of different curve segments in the connecting semi-algebraic path is at most  $\ell d^{O(k^2)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .  $\square$

A **uniform roadmap of**  $(Q, \mathcal{P})$  is a union of open curve segments and points satisfying the following two conditions:

- URM<sub>1</sub>: The signs of the polynomials  $P \in \mathcal{P}$  are constant on each curve segment,
- URM<sub>2</sub>: The intersection of this set with any basic closed semi-algebraic set

$$\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k)) = \{x \in \mathbb{R}^k \mid Q(x) = 0 \wedge \text{sign}(P(x)) \in \sigma(P)\},$$

where  $\sigma \in \{\{0\}, \{0, 1\}, \{0, -1\}\}^{\mathcal{P}}$  is a weak sign condition on  $\mathcal{P}$ , is a roadmap for  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$ .

As a first step we describe an algorithm which, given a polynomial  $Q$ , a set of polynomials  $\mathcal{P}$ , and a point  $p$  in  $\text{Zer}(Q, \mathbb{R}^k)$  constructs a finite number of continuous semi-algebraic curves starting at  $p$  so that every semi-algebraically connected component of every realizable sign condition of  $\mathcal{P}$  in  $\text{Zer}(Q, \mathbb{R}^k)$  sufficiently near and to the left of  $p$  contains one of these curves without the point  $p$ .

If  $p \in \text{Zer}(Q, \mathbb{R}^k)$ , denote by  $\text{SIGN}(\mathcal{P}, p)$  the set of sign conditions  $\sigma$  such that

$$\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k)) \cap B(p, r)_{<v}$$

is non-empty for all sufficiently small  $r > 0$ .

*Algorithm 16.3. [Linking Paths]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,
  - a finite set of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , in strong  $\ell$ -general position with respect to  $Q$ ,
  - a point  $p \in \text{Zer}(Q, \mathbb{R}^k)_z$ , described by a real univariate representation of degree at most  $d^{O(k)}$  with coefficients in  $D$ .

- **Output:** a finite set of semi-algebraic paths starting at  $p$  such that for some sufficiently small  $r$  and for every  $\sigma$  in  $\text{SIGN}(\mathcal{P}, p)$  every connected component of  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k)) \cap B(p, r)_{<v}$  contains one of these semi-algebraic paths (without the endpoint  $p$ ).
- **Complexity:**  $(s + 2^\ell) d^{O(k)}$ , where  $s$  is a bound on the number of polynomials in  $\mathcal{P}$  and  $d$  is a bound on the degree of  $Q$  and the polynomials in  $\mathcal{P}$ .
- **Procedure:**
  - Let  $\mathcal{P}(p)$  be the set of polynomials in  $\mathcal{P}$  (possibly empty) that are zero at  $p$  and let  $B(p, \varepsilon)$  be a ball of radius  $\varepsilon$  and center  $p$ , where  $\varepsilon$  is a new variable. Using Algorithm 13.1 (Sampling) with the polynomials defining  $B(p, \varepsilon)_{<v}$  along with the polynomials  $\mathcal{P}(p)$  as input and structure  $D[\varepsilon] \subset \mathbb{R}(\varepsilon)$ , find a finite set of points  $\mathcal{S}(\varepsilon)$  intersecting every semi-algebraically connected component of every realizable sign condition of the polynomials in  $\mathcal{P}$  in  $B(p, \varepsilon)_{<v}$ .
  - For every  $u = (f, g_0, g_1, \dots, g_k) \in \mathcal{S}(\varepsilon)$ , apply Algorithm 11.20 (Removal of Infinitesimals) with input  $f$  and  $\mathcal{P}(p)_u$ , output  $t(u)$  and define  $t_0 = \min_{u \in \mathcal{S}(\varepsilon)} t(u)$ . Replacing  $\varepsilon$  by  $t \in (0, t_0]$  defines for each  $u(\varepsilon) \in \mathcal{S}(\varepsilon)$ , with associated point  $q(\varepsilon)$ , a semi-algebraic path  $\gamma(u)$  such that  $\gamma(0) = p$ .

**Proof of correctness:** The semi-algebraic paths  $\gamma(u)$ ,  $u \in \mathcal{S}(\varepsilon)$ , join  $p$  to points in every semi-algebraically connected component of the realizable sign conditions of  $\mathcal{P}$  intersected with  $\text{Zer}(Q, \mathbb{R}^k) \cap B(p, \varepsilon)_{<v}$ . □

**Complexity analysis:** Since at most  $\ell$  polynomials can be zero at  $p$ , the complexity is  $(s + 2^\ell) d^{O(k)}$ , using Remark 13.3.

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ . □

We are going to describe now the Uniform Roadmap Algorithm. The algorithm will call itself recursively. In each recursive call the number of variables will strictly decrease. The base case when  $k = 1$  or  $\ell = 0$  are easy.

Note that if  $\ell = 0$ , then a roadmap  $\text{RM}(\text{Zer}(Q, \mathbb{R}^k))$  is a uniform roadmap for  $(Q, \mathcal{P})$  since on every semi-algebraically connected component of  $\text{Zer}(Q, \mathbb{R}^k)$  the signs of the polynomials in  $\mathcal{P}$  are fixed.

If  $k = 1$ , the zeroes of  $Q$  are isolated, the roadmap consists of the zeroes of  $Q$ .

*Algorithm 16.4.* **[Uniform Roadmap]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ ,

- a finite set  $\mathcal{P} \subset \mathbb{D}[X_1, \dots, X_k]$ ,
- a natural number  $\ell$  such that  $\mathcal{P}$  is in strong  $\ell$ -general position with respect to  $Q$ .
- **Output:** a semi-algebraic set  $\text{URM}(Q, \mathcal{P})$  satisfying conditions  $\text{URM}_1$  and  $\text{URM}_2$ . Moreover,  $\text{URM}(Q, \mathcal{P})$  is described by real univariate representations and curve segments representations, and each of these representations is labeled by a subset  $\mathcal{R} \subset \mathcal{P}$  such that its associated point or curve segment is contained in  $\text{Zer}(\mathcal{R} \cup \{Q\}, \mathbb{R}^k)$ .
- **Complexity:**  $s^{\ell+1} d^{O(k)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Initialize  $i := 1, \mathcal{T}, \sigma := \emptyset, S := Q^2, \mathcal{R} := \mathcal{P}, m := \ell$ .
  - Step 1: ( $\star$ ) For each  $\mathcal{R}' \subset \mathcal{R}, \#(\mathcal{R}') \leq m$ .
    - Step 1 a): If  $\#(\mathcal{R}') = m$ , describe the isolated zeroes of  $\mathcal{T}, S + \sum_{P \in \mathcal{R}'} P^2$ . These points are placed in a set of distinguished points for  $\mathcal{R}'$ . A distinguished value for  $\mathcal{R}'$  is the  $i$ -th coordinate of a distinguished point for  $\mathcal{R}'$ .
    - Step 1 b): If  $\#(\mathcal{R}') < m$ , run Algorithm 15.2 (Curve Segments) with  $\mathcal{T}, \sigma, S + \sum_{P \in \mathcal{Q}'} P^2$ , and an empty set of univariate representations as input.

Label each curve segment by  $\mathcal{R}'$ . The endpoints of these curve segments are labeled by  $\mathcal{R}'$  and are placed in a set of distinguished points for  $\mathcal{R}'$ . A distinguished value for  $\mathcal{R}'$  is the  $i$ -th coordinate of a distinguished point for  $\mathcal{R}'$ .

- Step 1 c): Run Algorithm 16.1 (Special Values) for  $\mathcal{T}, \sigma, S + \sum_{P \in \mathcal{R}'} P^2$  and intersect the curve segments obtained in Step 1 a) with the corresponding special hyperplanes. Add these points to the set of distinguished points for  $\mathcal{R}'$ . Append their  $i$ -th coordinates to the set of distinguished values for  $\mathcal{R}'$ .
- Step 1 d): Compute the intersection of each curve segment output in Step 1 a) with  $\text{Zer}(P, \mathbb{R}^k)$  for each  $P \in \mathcal{R}$ . Note that the intersection of a curve segment with the zero set of a polynomial, is either the segment itself, or a finite set of points (possibly empty). This is checked by substituting the parametrized univariate representation of the curve into each polynomial in  $\mathcal{R}$  and checking whether the resulting univariate polynomial vanishes identically or not.

If the intersection is the curve segment itself, ignore this intersection. Otherwise, the points of intersection yield a partition of the curve segment. Add these points to the set of distinguished points for  $\mathcal{R}' \cup \{P\}$ . Append their  $i$ -th coordinates to the set of distinguished values for  $\mathcal{R}' \cup \{P\}$ . Store the sign vector of the set of polynomials  $\mathcal{R}$  on each curve segment and point computed above.

- Step 2: For every distinguished point  $p$  with label  $\mathcal{R}'$  and  $i$ -th coordinate  $v$ , add to the distinguished points for  $\mathcal{R}'$  the intersections of the hyperplane  $H_v$  with the curves constructed for  $\mathcal{R}''$  in Step 1, where  $\mathcal{R}' \subset \mathcal{R}'' \subset \mathcal{R}$ ,  $\#(\mathcal{R}'') \leq m - \#(\mathcal{R}') - 1$ .
- Step 3: For all distinguished value  $v$  specified by  $A, \alpha$ , call the algorithm recursively, returning to  $(\star)$  with input

$$\begin{aligned} i &:= i + 1, \\ \mathcal{T}, \sigma &:= \mathcal{T}, A, \sigma, \alpha, \\ S &:= S + \sum_{P \in \mathcal{P}'} P^2, \\ \mathcal{R} &:= \mathcal{R} \setminus \mathcal{R}', \\ m &:= m - \#(\mathcal{R}'). \end{aligned}$$

Denote by  $\text{URM}_0(Q, \mathcal{P})$  the output so obtained.

- Step 4: For each distinguished point  $p$  and the corresponding distinguished hyperplane, use Algorithm 16.3 (Linking Paths) to construct semi-algebraic paths joining  $p$  to points in every semi-algebraically connected component of every realizable sign condition of the set of polynomials  $\mathcal{P}$  intersected with  $\text{Zer}(Q, \mathbb{R}^k) \cap B(p, r)_{<\pi(p)}$ , for some small enough  $r$ . Let the other endpoints of these curves be a finite set  $\mathcal{S}$ . Connect the points of  $\mathcal{S}$  to some  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$  using Algorithm 16.2 (Bounded Connecting).
- Output all the curve segments and distinguished points, each labeled by the sign condition it satisfies. This is the set  $\text{URM}(Q, \mathcal{P})$ .

The proof of correctness of Algorithm 16.4 (Uniform Roadmap) is based on the following results.

Let  $S$  be the semi-algebraic set defined by  $Q = 0, P \geq 0, P \in \mathcal{P}$ , and let  $\text{RM}(S) = S \cap \text{URM}(Q, \mathcal{P})$ .

**Proposition 16.7.** *The set  $\text{RM}(S)$  is a roadmap for the set  $S$ .*

**Proof:** We first show that  $\text{RM}(S)$  satisfies  $\text{RM}_2$ .

For any  $x \in \mathbb{R}$  such that  $S_x$  is non-empty, and for any semi-algebraically connected component  $C$  of  $S_x$ , there exists a semi-algebraically connected component  $C'$  of a non-empty algebraic set,  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)_x$  such that  $C' \subset C$  (see Proposition 13.1). Since, in Step 1 b of the algorithm we construct curves using the Algorithm 15.2 (Curve Segments) on all non-empty algebraic sets of the form  $\text{Zer}(\mathcal{P}', \mathbb{R}^k)$ , it is clear that  $\text{RM}(S)$  intersects  $C$ . Thus  $\text{RM}(S)$  satisfies  $\text{RM}_2$ .

We next show that  $\text{RM}(S)$  satisfies condition  $\text{RM}_1$  as well. This is the content of the following two lemmas. Let  $v(1), \dots, v(\ell)$  be the set of distinguished values computed by the algorithm.



**Lemma 16.8.** *For  $1 \leq i \leq \ell$ , if  $\text{RM}(S)_{<v(i)}$  satisfies condition  $\text{RM}_1$  for the set  $S_{\leq v(i)}$  then,  $\text{RM}(S)_{<v(i+1)}$  satisfies condition  $\text{RM}_1$  for the set  $S_{<v(i+1)}$ .*

**Proof:** Let  $C$  be a semi-algebraically connected component of  $S_{<v(i+1)}$  and let  $\Gamma$  be a semi-algebraically connected component of  $\text{RM}(S) \cap C_{[v(i), v(i+1))}$ . The set  $\Gamma_{v(i)}$  is non-empty since there is no distinguished value in  $(v(i), v(i+1))$ . It is then clear that  $\text{RM}(S) \cap C_{\leq v(i)} \cup \Gamma$  is semi-algebraically connected. Since  $\text{RM}(S) \cap C_{\leq v(i)}$  is semi-algebraically connected, the conclusion follows.  $\square$

**Lemma 16.9.** *For  $1 \leq i \leq \ell$ , if  $\text{RM}(S)_{<v(i)}$  satisfies condition  $\text{RM}_1$  for the set  $S_{<v(i)}$ , then  $\text{RM}(S)_{\leq v(i)}$  satisfies condition  $\text{RM}_1$  for the set  $S_{\leq v(i)}$ .*

**Proof:** Let  $C$  be a semi-algebraically connected component of  $S_{\leq v(i)}$ . We prove that  $\text{RM}(S) \cap C$  is semi-algebraically connected.

Let  $B_1, \dots, B_h$  be the semi-algebraically connected components of  $S \cap C_{<v(i)}$ . Then, by Lemma 16.3,  $C = C_1 \cup C_2 \dots \cup C_N$ , where each  $C_i$  is either  $B_j$  or a semi-algebraically connected component of  $\text{Zer}(\mathcal{P}', \mathbf{R}^k) \cap S_{v(i)}$ , for some  $I \subset \{1, \dots, s\}$ , where  $v(i)$  is an  $X_1$ -special value of  $\text{Zer}(\mathcal{P}', \mathbf{R}^k)$ .

Let  $\Gamma = \text{RM}(S) \cap C$  and  $\Gamma(i) = \text{RM}(S) \cap C_i$  for  $1 \leq i \leq N$ . Then  $\Gamma = \bigcup_i \Gamma(i)$ .

First, we claim that each  $\Gamma(j)$  is semi-algebraically connected. If  $C_j$  is a semi-algebraically connected component of  $\text{Zer}(\mathcal{P}') \cap S_{v(i)}$  for some  $I \subset \{1, \dots, s\}$  containing  $Q$ , then, since  $v(i)$  is an  $X_1$ -special value for this algebraic set,  $\Gamma(j)$  is semi-algebraically connected by Step 4 of the algorithm.

Else, by the hypothesis of the lemma, we know that  $\Gamma(j)_{<v(i)}$  is semi-algebraically connected. Thus,  $\Gamma(j)$  can have at most one semi-algebraically connected component whose intersection with  $(\mathbf{R}^k)_{<v(i)}$  is non-empty, and all the other semi-algebraically connected components of  $\Gamma(j)$  must lie in  $\pi^{-1}(v(i))$ . Hence each of these must contain a distinguished point. But, by Step 4 of the algorithm, the distinguished points get connected to  $\Gamma(j)_{<v(i)}$ . Thus,  $\Gamma(j)$  can have only one semi-algebraically connected component.

Moreover, if  $C_j \cap C'_j \neq \emptyset$ , then  $\Gamma(j)$  and  $\Gamma(j')$  are connected in  $\text{RM}(S)$ . This is so since, according to Lemma 16.4,  $C_j \cap C'_j$  intersects an algebraic set which has  $v(i)$  as an  $X_1$ -pseudo-critical value and thus contains a distinguished point which gets linked to both  $\Gamma(j)$  and  $\Gamma(j')$ .

It follows that  $\Gamma$  is semi-algebraically connected. This proves the lemma.  $\square$

The proposition now follows by induction on  $i$ .  $\square$

**Proof of correctness:** Note that Step 2 b and Step 3 of the algorithm make it evident that  $\text{RM}(Q, \mathcal{P})$  satisfies condition  $\text{URM}_1$ . That it satisfies condition  $\text{URM}_2$  follows from Proposition 16.7.  $\square$

**Complexity analysis:** When  $\ell = 1$ , it follows from the analysis of the algebraic case that the number of arithmetic operations is  $s d^{O(k^2)}$ . When  $k = 1$ , the number of arithmetic operations is  $s d^{O(1)}$ .

In Step 1 a) the total number of arithmetic operations in D is  $s \binom{s}{\ell} d^{O(k)} = s^{\ell+1} d^{O(k)}$ .

In Step 1 b, the total number of calls to Algorithm 15.2 (Curve Segments) is  $\sum_{j=1}^{\ell-1} \binom{s}{j}$ , and each call costs  $d^{O(k)}$ . Thus, the total cost of the calls to Algorithm 15.2 (Curve Segments) is bounded by  $\sum_{j=1}^{\ell-1} \binom{s}{j} d^{O(k)}$  arithmetic operations in D.

In Step 1 c, the cost of each call to Algorithm 16.1 (Special Values) as well as the cost of computing the intersection of each curve segment with the special hyperplanes are bounded by  $d^{O(k)}$ , and hence the total cost of this step is bounded by  $\sum_{j=1}^{\ell-1} \binom{s}{j} d^{O(k)}$ .

In Step 1d, the cost of computing the intersection of each curve segment computed with the zero sets of each polynomial in  $\mathcal{P}$  is bounded by  $s d^{O(k)}$ , and the total cost of Step 1 d, for all  $I$  considered, is  $s \sum_{j=1}^{\ell-1} \binom{s}{j} d^{O(k)}$ .

In Step 2, the cost is  $\sum_{j=1}^{\ell-1} \binom{s}{j} 2^j d^{O(k)}$ , since  $\sum_{j=1}^{\ell-1} \binom{s}{j} 2^j$  is the number of pairs  $(\mathcal{R}', \mathcal{R}'')$  considered.

Note that the combinatorial level in the recursive call is at most  $\ell - \#(\mathcal{P}') - 1$  and the number of variables is  $k - 1$ .

We now count the recursive calls. For each  $j$ ,  $0 \leq j \leq \ell - 1$ , we make  $\sum_{j=1}^{\ell-1} \binom{s}{j} d^{O(k)}$  recursive calls to the algorithm with combinatorial level  $\ell - j$  and ambient space dimension  $k - 1$ .

Let  $T(s, d, \ell, k)$  denote the complexity of the algorithm with these parameters. Since at any depth of the recursion the cost of a single arithmetic operations is bounded by  $d^{O(k^2)}$  arithmetic operations in D, we ignore the fact that the ring changes as we go down in the recursion. Thus, we have the following recurrence,

$$\begin{aligned}
 T(s, d, \ell, k) &\leq \sum_{j=1}^{\ell-1} \binom{s}{j} d^{O(k)} T(s - j, d, \ell - j, k - 1) + s \sum_{j=0}^{\ell} \binom{s}{j} d^{O(k)}, \\
 &\qquad \qquad \qquad \ell > 0, k > 1, \\
 T(s, d, 0, k) &= s d^{O(k^2)}, \quad k > 1, \\
 T(s, d, \ell, 1) &= s d^{O(1)}.
 \end{aligned}$$

This recurrence solves to  $T(s, d, \ell, k) = s^{\ell+1} d^{O(k^2)}$ .

In Step 4 the total cost of the calls to the Algorithm 16.3 (Linking Paths) and Algorithm 16.2 (Bounded Connecting) is bounded by  $s^{\ell+1} d^{O(k)}$ .

It follows immediately that the total cost is still bounded by  $s^{\ell+1} d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

### 16.3 Computing Connected Components of Sign Conditions

For complexity reasons, the formulas describing the semi-algebraically connected components of a given semi-algebraic set produced by our algorithm will not necessarily be written as disjunctions of sign conditions. This differs from some of our previous algorithms (such as eliminating quantifiers, or describing the semi-algebraically connected components of algebraic sets).

**Notation 16.10.** Let  $\phi(Y)$  be a quantifier free formula,

$$T, \sigma, u = (T_\ell, g_0, g_\ell, \dots, g_k)$$

a parametrized real univariate triangular representation with parameters  $Y = (Y_1, \dots, Y_k)$ . With  $T = (T_1, \dots, T_\ell)$ , we denote by  $\phi_u(Y)$  the formula

$$(\forall T) \left( \bigwedge_{1 \leq i \leq \ell} \left( T_i(Y, T_1, \dots, T_i) = 0 \wedge \bigwedge_{h \in \text{Der}(T)} \text{sign}(h(Y, T)) = \sigma(h) \right) \right) \Rightarrow \phi(u)$$

where  $\phi(u)$  is obtained by replacing  $Y_j$  by  $g_j(Y, T_1, \dots, T_\ell)/g_0(Y, T_1, \dots, T_\ell)$ ,  $j \geq \ell$  and clearing denominators. □

*Algorithm 16.5. [Parametrized Bounded Connecting]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$ , such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$ ,
  - a finite set of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$  in strong  $k$ -general position with respect to  $Q$ .
- **Output:**
  - a finite set of polynomials  $\mathcal{A}$  containing  $\mathcal{P}$ ,
  - a finite set  $\Theta$  of  $\mathcal{A}$ -quantifier free formulas such that for every semi-algebraically connected component  $S$  of the realization of every weak sign condition on  $\mathcal{P}$  on  $\text{Zer}(Q, R^k)$ , there exists a subset  $\Theta(S) \subset \Theta$  such that

$$S = \bigcup_{\theta \in \Theta(S)} \text{Reali}(\theta, \text{Zer}(Q, R^k)),$$

- for every  $\theta \in \Theta$ , a parametrized path  $\Gamma(\theta) \subset R^{2k}$  such that  $\Gamma(\theta)_y$  is a semi-algebraic set of dimension at most one, that connects for every  $y \in \text{Reali}(\theta)$  the point  $y$  to some roadmap  $\text{RM}(\text{Zer}(\mathcal{P}' \cup \{Q\}, R^k))$  where  $\mathcal{P}' \subset \mathcal{P}$ , staying inside

$$\text{Reali}(\bar{\sigma}(y), \text{Zer}(Q, R^k)).$$

Moreover, for every  $y \in \text{Reali}(\theta, \text{Zer}(Q, R^k))$ , the description of  $\Gamma(\theta)_y$  is fixed and the endpoint of  $\Gamma(\theta)_y$ , described by the real univariate representation  $w(\theta)$  is independent of  $y$ .

- **Complexity:**  $s^{\ell+1} d^{O(k^4)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**

- Step 1: Fix an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$ , such that  $\#(I) \leq \ell$ , and denote by  $\mathcal{P}_I$  the set of polynomials  $\{Q\} \cup \{P_i \in \mathcal{P} \mid i \in I\}$ . If  $\text{Zer}(\mathcal{P}_I, \mathbb{R}^k) \neq \emptyset$ , compute using Algorithm 15.12 (Parametrized Bounded Algebraic Connecting) a family of polynomials  $\mathcal{A}(I) \subset D[Y_1, \dots, Y_k]$ , and for every  $\rho \in \text{SIGN}(\mathcal{A}(I), \mathcal{P}_I)$  a semi-algebraic set  $\Gamma(\rho) \subset \mathbb{R}^{2k}$  such that, for every  $y \in \text{Reali}(\rho)$ ,  $\Gamma(\rho)_y$  connects the point  $y$  to a distinguished point of  $\text{RM}(\text{Zer}(\mathcal{P}_I, \mathbb{R}^k))$ , described by the real univariate representation  $w(\rho)$ . Moreover, for every  $y \in \text{Reali}(\rho)$ , the description of  $\Gamma(\rho)_y$  is fixed.
- Step 2: Fix an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$ , such that  $\#(I) \leq \ell$ , and  $j \in \{1, \dots, s\} \setminus I$ . Compute a family of polynomials whose signs control the manner in which  $\text{Zer}(P_j, \mathbb{R}^k)$  intersects  $\Gamma(\rho)$ . More precisely, compute a family of polynomials  $\mathcal{B}(\rho, j)$  containing  $\mathcal{A}(I)$  and the subset  $\Sigma(\rho, j)$  of elements  $\rho' \in \text{SIGN}(\rho, \mathcal{B}(\rho, j))$  such that for every  $y \in \text{Reali}(\rho')$ ,
  - the intersection of  $\Gamma(\rho)(y)$  with  $\text{Zer}(P_j, \mathbb{R}^k)$  is non-empty,
  - the Thom encodings describing the various points of intersection of  $\Gamma(\rho)_y$  with  $\text{Zer}(P_j, \mathbb{R}^k)$  remain constant.

In order to achieve this, we first use Algorithm 14.7 (Parametrized Triangular Thom Encoding) as follows. Each curve segment of  $\Gamma(\rho)$  is described by:

- a parametrized triangular Thom encoding  $\mathcal{T}(Y, X_{\leq i}), \sigma$ ,
- a parametrized univariate representation with parameters  $(X_{\leq i})$ ,
 
$$u = (f(Y, X_{\leq i}, T), g_0(Y, X_{\leq i}, T), g_{i+1}(Y, X_{\leq i}, T), \dots, g_k(Y, X_{\leq i}, T)),$$
- a sign condition on  $\text{Der}(f)$ .

For each such curve segment in  $\Gamma(\rho)$ , first compute

$$\text{RElim}_T(P_{j,u}, f) \subset D[Y, X_{\leq i}]$$

using Algorithm 11.19. Then call Algorithm 14.7 (Parametrized Triangular Thom Encoding) with input  $\mathcal{T} \cup \{P\}$  for each  $P \in \text{RElim}_T(P_{j,u}, f)$ .

The output is:

- a finite set  $\mathcal{B}' \subset D[Y]$ ,
- for every  $\rho' \in \text{SIGN}(\mathcal{B}')$ , a list of sign conditions on  $\text{Der}(\mathcal{T})$  specifying, for every  $y \in \text{Reali}(\rho')$ , the list of triangular Thom encodings of the roots of  $\mathcal{T}(y) \cup \{P(y)\}$ .

Let  $\mathcal{B}(\rho, j)$  be the union of all the  $\mathcal{B}'$  obtained above along with  $\mathcal{A}(I)$ . Now use Algorithm 15.7 (Parametrized Comparison of Roots) to the order the various points of intersections and compute  $\Sigma(\rho, j)$ .

- Step 3: Fixing an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$ , such that  $\#(I) \leq \ell$ , denote by  $\Phi(I)$  the set of formulas

$$\mathcal{F}(\rho) := \left( Q = 0 \wedge \bigwedge_{i \in I} P_i(x) = 0 \wedge \bigwedge_{A \in \mathcal{A}(I)} \text{sign}(A)(x) = \rho(A) \right)$$

for all  $\rho \in \text{SIGN}(\mathcal{A}(I))$ . Similarly, fixing an ordered tuple of indices  $I$  with element in  $\{1, \dots, s\}$ , with  $\#(I) \leq \ell - 1$ , and  $j \in \{1, \dots, s\} \setminus I$ , denote by  $\Phi(I, j)$  the set of formulas

$$\mathcal{F}(\rho) \wedge \bigwedge_{B \in \mathcal{B}(\rho, j)} \text{sign}(B)(x) = \rho'(B),$$

for all  $\rho \in \text{SIGN}(\mathcal{A}(I))$  and  $\rho' \in \Sigma(\rho, j)$ .

For every  $\phi \in \Phi(I, j)$  and every  $y \in \text{Reali}(\phi)$ , the first point of intersection of  $\Gamma(\phi)_y$  with  $\text{Zer}(P_j, \mathbb{R}^k)$ ,  $F(\phi)(y)$ , is described by a real parametrized univariate representation with parameters  $Y$ , denoted by  $u(\phi)(Y), \tau(\phi)$ .

Denote by  $\gamma(\phi)_y$  the part of the semi-algebraic path  $\Gamma(\phi)_y$ , starting at  $y$  and ending at  $F(\phi)(y)$ , and by  $\gamma(\phi) \subset \mathbb{R}^{2k}$  the union of  $\{y\} \times \gamma(\phi)_y$  for  $y \in \text{Reali}(\phi)$ .

Compose the functions  $F(\phi)$  inductively, as follows. Fix an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$ , such that  $\#(I) \leq \ell$  and  $\text{Zer}(\mathcal{P}_I, \mathbb{R}^k) \neq \emptyset$  and initialize  $\Psi(I) := \Phi(I)$  and for every  $\psi \in \Phi(I)$  associated to  $\rho$ ,  $v(\psi)(Y) := Y$ ,  $w(\psi) := w(\rho)$ ,  $\Gamma(\psi) = \emptyset$ .

Fix an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$ , such that  $\#(I) \leq \ell - 1$ . We will compute for every  $J$ ,  $1 \leq \#(J) \leq \ell - \#(I)$  a finite set of quantifier free formulas  $\Psi(I, J)$  and for every  $\psi \in \Psi(I, J)$ , a parametrized real univariate triangular representation  $\mathcal{T}(\psi)$ ,  $\sigma(\psi)$ , as well as  $v(\psi)$ ,  $w(\psi)$ ,  $\Gamma(\psi)$ .

Let  $J$  be an ordered tuple of indices with elements in  $\{1, \dots, s\}$  such that  $1 \leq \#(J) \leq \ell - \#(I)$  and suppose that a finite set of quantifier free formulas  $\Psi(I, J)$ , as well as for every  $\psi \in \Psi(I, J)$ , parametrized real univariate triangular representation  $\mathcal{T}(\psi), \sigma(\psi), v(\psi), w(\psi), \Gamma(\psi)$ , have already been computed.

Let  $\bar{J} = J \cdot j$ ,  $j \in \{1, \dots, s\} \setminus I \cdot J$ , and define the set  $\Psi(I, \bar{J})$  as follows.

For each  $\psi \in \Psi(I, J)$ , each  $\phi_1 \in \Phi(I \cdot J, j)$ , and each  $\phi_2 \in \Phi(I \cdot J \cdot j)$ , let  $v = u(\phi_1)_{v(\psi)}$ . Compute, using Algorithm 14.5 (Quantifier Elimination), a quantifier-free formula  $\overline{\phi_{1,v(\psi)} \wedge \phi_{2,v}}$  equivalent to  $\phi_{1,v(\psi)} \wedge \phi_{2,v}$ . Include in  $\Psi(I, \bar{J})$  all  $\psi \wedge \overline{\phi_{1,v(\psi)} \wedge \phi_{2,v}}$  which are realizable using Algorithm 13.1 (Computing Realizable Sign Conditions), with input the family of polynomials appearing in  $\psi \wedge \overline{\phi_{1,v(\psi)} \wedge \phi_{2,v}}$ .

For every  $\psi' = \psi \wedge \overline{\phi_{1,v(\psi)} \wedge \phi_{2,v}} \in \Psi(I, \bar{J})$ , define  $v(\psi')$  as  $v = (\mathcal{T}_{i+1}, \bar{g}_0, \dots, \bar{g}_k)$ , and define a new triangular system  $\mathcal{T}(\psi')$  obtained by appending  $\mathcal{T}_{i+1}$  to  $\mathcal{T}(\psi)$  and a new list of sign vectors  $\sigma(\psi')$  by appending  $\tau(\phi_1)$  to the list  $\sigma(\psi)$ . Finally, let  $w(\psi') := w(\phi_2)$ .

Define  $\Gamma(\psi') = \Gamma(\psi) \cup \gamma(\phi_{1,v(\psi)})$ .

- Step 4: For an ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$  such that  $\#(I) \leq \ell - 1$ , and an ordered tuple of indices  $J$  such that  $1 \leq \#(J) \leq \ell - \#(I)$  with elements in  $\{1, \dots, s\}$  and a formula  $\psi \in \Psi(I, J)$  the semi-algebraic path  $\Gamma(\psi) \cup \Gamma(\rho)(v(\psi))$ , where  $\rho$  is the sign condition on  $\mathcal{A}_{I,J}$  satisfied at  $v(\psi)$ , may or may not be a valid connecting semi-algebraic path depending on whether any polynomials in  $\mathcal{P} \setminus \mathcal{P}_{I,J}$  vanish on any one of its segments.

Compute the formula  $\theta(\psi, j)$  expressing the conditions on  $Y$  ensuring that  $P_j$  does not vanish on  $\Gamma(\psi) \cup \Gamma(\rho)(v(\psi))$ , using Algorithms 14.7 (Parametrized Triangular Thom Encoding) and 15.7 (Parametrized Comparison of Roots) for all the real parametrized univariate representations describing  $\Gamma(\psi) \cup \Gamma(\rho)(v(\psi))$ .

Define the set  $\Theta(I, J)$  of formulas  $\psi \wedge \bigwedge_{j \notin I,J} \theta(\psi, j)$  with  $\psi \in \Psi(I, J)$ , and, for  $\theta \in \Theta(I, J)$ ,

$$w(\theta) = w(\psi), \Gamma(\theta) = \Gamma(\psi) \cup \Gamma(\rho)(v(\psi))$$

$w(\theta) = w(\psi)$ ,  $\Gamma(\theta) = \Gamma(\psi) \cup \Gamma(\rho)(v(\psi))$ , where  $\rho$  is the sign condition on  $\mathcal{A}_{I,J}$  satisfied at  $v(\psi)$ .

Since formulas of  $\Theta(I, J)$  are refinements of formulas in  $\Psi(I, J)$ , every  $\theta \in \Theta(I, J)$  defines a subset  $\Gamma(\theta)$  such that  $\Gamma(\theta)(y)$  is a path connecting  $y$  to the point of the roadmap for  $\mathcal{P}_{I,J}$  inside  $\text{Reali}(\sigma(y), \text{Zer}(Q, \mathbb{R}^k))$ , described by the real univariate representation  $w(\theta)$ .

Define  $\Theta$  as the union for every ordered tuple of indices  $I$  with elements in  $\{1, \dots, s\}$  such that  $\#(I) \leq \ell - 1$ , and every ordered tuple of indices  $J$  with elements in  $\{1, \dots, s\} \setminus I$  such that  $1 \leq \#(J) \leq \ell - \#(I)$  of  $\Theta(I, J)$ .

Define  $\mathcal{A} \subset \mathbb{D}[X_1, \dots, X_k]$  to be the set of polynomials appearing in the formulas of  $\Theta$ .

**Proof of correctness :** It is clear from the algorithm that each formula  $\theta$  obtained in Step 4 has the property that every  $y \in \text{Reali}(\theta)$  gets connected to a unique distinguished point of some algebraic set  $\text{Zer}(\mathcal{P}_I)$  by  $\Gamma(\theta)_y$  inside  $\text{Reali}(\bar{\sigma}(y), \text{Zer}(Q, \mathbb{R}^k))$ . Thus, each  $\text{Reali}(\theta)$  must be fully contained in some connected component of a realizable weak sign condition of  $\mathcal{P}$ . Moreover, clearly  $\bigcup_{\theta \in \Theta} \text{Reali}(\theta) = \text{Zer}(Q, \mathbb{R}^k)$ . It is easy to see that  $\Gamma(\theta)$  is a parametrized path, by the correctness of Algorithm 15.12 (Parametrized Bounded Algebraic Connecting).  $\square$

**Complexity analysis:** Using the complexity analysis of Algorithm 15.12 the complexity of Step 1 is bounded by  $\sum_{i=1}^{\ell} \binom{s}{i} d^{O(k^3)}$ . The number and degrees of the polynomials in the various  $\mathcal{A}(I)$  are bounded by  $d^{O(k^2)}$ . Note that the number of elements of  $\Phi(I)$  coincides with the number of non-empty sign conditions on  $\mathcal{A}(I)$  and is bounded by  $d^{O(k^3)}$ .

Similarly, using the complexity analysis of Algorithms 14.7 and 15.7, the complexity of Step 2 is bounded by  $\sum_{i=1}^{\ell-1} (s-i) \binom{s}{i} d^{O(k^3)}$ . The number and degrees of the polynomials in the various  $\mathcal{B}(\rho, j)$  are bounded by  $d^{O(k^2)}$ . Note that the number of elements of  $\Phi(I, j)$  coincides with the number of non-empty sign conditions on  $\mathcal{B}(I, j)$  and is bounded by  $d^{O(k^3)}$ .

In Step 3 the complexity for an ordered tuple  $I$  of indices of length  $p$  is  $\sum_{i=1}^{\ell-p} \binom{s}{i} d^{O(k^4)}$  since there are  $\sum_{i=1}^{\ell-p} \binom{s}{i}$  choices for  $J$ . The degrees of the polynomials appearing in the formulas of  $\Psi(I, J)$  are bounded by  $d^{O(k^3)}$ . Thus, the total complexity of Steps 3 is bounded by  $s^\ell d^{O(k^4)}$ . The number of elements of  $\Theta$  is  $s^\ell d^{O(k^3)}$ . Moreover, for every  $\theta \in \Theta$ , since  $w(\theta)$ 's is a distinguished point of the roadmap of some  $\text{Zer}(P_I, \mathbb{R}^k)$ , the triangular system defining  $w(\theta)$  has polynomials of degree at most  $d^{O(k)}$ .

Finally, the complexity of Step 4 is bounded by  $s^{\ell+1} d^{O(k^4)}$ .

The complexity of the algorithm is  $s^{\ell+1} d^{O(k^4)}$ . The number of polynomials in the family  $\mathcal{A}$  is also  $s^{\ell+1} d^{O(k^4)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^4)}$ .  $\square$

### Algorithm 16.6. [Basic Connected Components]

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$  of polynomials of degree at most  $d$ .
- **Output:** quantifier free formulas whose realizations are the semi-algebraically connected components of  $\text{Reali}(\sigma, \mathbb{R}^k)$ , for the realizable sign conditions  $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$ .
- **Complexity:**  $s^{k+1} d^{O(k^4)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Take  $Q = \varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1$ .
  - Replace the set  $\mathcal{P}$  by the family  $\mathcal{P}^*$  defined by

$$\begin{aligned} P_i^* &= \{(1-\delta)P_i + \delta H_k(d', i), (1-\delta)P_i - \delta H_k(d', i), \\ &\quad (1-\delta)P_i + \delta \gamma H_k(d', i), (1-\delta)P_i - \delta \gamma H_k(d', i)\} \\ \mathcal{P}^* &= \{P_1^*, \dots, P_s^*\} \end{aligned}$$

for  $0 \leq i \leq s$ , where  $H_k(d', i) = (1 + \sum_{j=1}^k i^j X_j^{d'})$  and  $d' > d$ .

- For every non-empty sign condition  $\sigma$  on  $\mathcal{P}$

$$\begin{aligned} P_i &= 0, & i \in I \subset \{1, \dots, s\} \\ P_i &> 0, & i \in J \subset \{1, \dots, s\} \setminus I \\ P_i &< 0, & i \in \{1, \dots, s\} \setminus (I \cup J), \end{aligned}$$

let  $\sigma^*$  be the weak sign condition on  $Q$  and  $\mathcal{P}^*$  defined by

$$\begin{aligned} Q &= 0. \\ -\gamma \delta H_k(d', i) &\leq (1 - \delta) P_i \leq \gamma \delta H_k(d', i), \quad i \in I, \\ (1 - \delta) P_i &\geq \delta H_k(d', i), \quad i \in J \\ (1 - \delta) P_i &\leq -\delta H_k(d', i), \quad i \in \{1, \dots, s\} \setminus (I \cup J). \end{aligned}$$

Apply Algorithm 16.5 (Parametrized Bounded Connecting) with input  $Q$  and  $\mathcal{P}^*$  and output  $\Theta$ . Compute the set  $\Theta_\sigma$  of  $\theta \in \Theta$  such that  $w(\theta)$  belongs to the realization of  $\sigma^*$ . Using Algorithm 16.4 (Uniform Roadmap) with input  $Q$ ,  $\mathcal{P}^*$  and the  $w(\theta)$ ,  $\theta \in \Theta_\sigma$ , partition  $\{w(\theta) \mid \theta \in \Theta_\sigma\}$  into subsets  $W_1, \dots, W_r$  such that all points of  $W_i$  belong to the same semi-algebraically connected component of the realization of  $\sigma^*$ . Compute  $\Phi_1, \dots, \Phi_r$  with  $\Phi_i = \bigvee_{\{\theta \in \Theta_\sigma \mid w(\theta) \in W_i\}} \theta$ .

- For every semi-algebraically connected component  $C$  of the realization of  $\sigma$  there exists a semi-algebraically connected component  $C'$  of the realization of  $\sigma^*$  such that  $\pi(C') \cap \mathbb{R}^k = C$ . Consider the formula  $\Phi_i(X_1, \dots, X_k, X_{k+1})$  describing  $C'$  and denote by  $\Psi_i(X_1, \dots, X_k)$  the formula obtained by replacing each atom  $F(X_1, \dots, X_{k+1})$  of  $\Phi_i$  by the quantifier free formula equivalent to  $(\exists X_{k+1}) X_{k+1} < 0 \wedge F(X_1, \dots, X_{k+1})$  using Algorithm 14.5 (Quantifier Elimination). The formula  $\Psi_i$  describes  $\pi(C')$ . Then  $\text{Remo}_{\varepsilon, \delta, \gamma}(\Psi_i(Y))$  (Notation 14.6) defines  $C$ .

**Proof of correctness:** According to Proposition 13.7, every semi-algebraically connected component  $C$  of every strict sign condition of the original family  $\mathcal{P}$  corresponds to a semi-algebraically connected component  $C'$  of a weak sign condition on  $\mathcal{P}^*$ . Moreover,  $C = \pi(C') \cap \mathbb{R}^k$ . Now use Proposition 14.7.  $\square$

**Complexity analysis:** The family  $\mathcal{P}^*$  has combinatorial level  $k$  by Proposition 13.6. Using the complexity analysis of Algorithm 16.5 (Parametrized Bounded Connecting), the complexity of computing  $\Theta$  is  $s^{k+1}d^{O(k^2)}$ . Applying Algorithm 16.4 costs  $s^{k+1}d^{O(k^2)}$ , and the points  $w(\theta)$ ,  $\theta \in \Theta$  are distinguished points of this uniform roadmap. For every atom, the quantifier elimination performed costs  $d^{O(k^4)}$ , since there is one variable to eliminate,  $k$  free variables and one polynomial of degree  $d^{O(k^3)}$ , according to the complexity of Algorithm 14.5 (Quantifier Elimination). The total number of the atoms to consider is  $s^k d^{O(k^4)}$ .



So the total complexity is  $s^{k+1}d^{O(k^4)}$ . The degrees of the polynomials that appear in the output are bounded by  $d^{O(k^3)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^4)}$ .  $\square$

**Theorem 16.11.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset D[X_1, \dots, X_k]$  with  $\deg(P_i) \leq d$ , for  $1 \leq i \leq s$ . There exists an algorithm that outputs quantifier-free semi-algebraic descriptions of all the semi-algebraically connected components of every realizable sign condition of the family  $\mathcal{P}$ . The complexity of the algorithm is bounded by  $s^{k+1}d^{O(k^4)}$ . The degrees of the polynomials that appear in the output are bounded by  $d^{O(k^3)}$ . Moreover, if the input polynomials have integer coefficients whose bitsize is bounded by  $\tau$  the bitsize of coefficients output is  $\tau d^{O(k^3)}$ .*

## 16.4 Computing Connected Components of a Semi-algebraic Set

We first construct data for adjacencies for  $\mathcal{P}$  on  $\text{Zer}(Q, \mathbb{R}^k)$ , ensuring that if the union of two semi-algebraically connected components of two different sign conditions for  $\mathcal{P}$  on  $\text{Zer}(Q, \mathbb{R}^k)$  is semi-algebraically connected, a path starting in a sign condition and ending in the other is constructed.

A set  $N$  of **data for adjacencies** for  $\mathcal{P}$  on  $\text{Zer}(Q, \mathbb{R}^k)$  is a set of triples  $(p, q, \gamma)$ , where  $p, q \in \text{Zer}(Q, \mathbb{R}^k)$ , and  $\gamma$  is semi-algebraic path joining  $p$  to  $q$  inside  $\text{Zer}(Q, \mathbb{R}^k)$ , such that for any two semi-algebraically connected components,  $C$  and  $D$  of  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$  and  $\text{Reali}(\tau, \text{Zer}(Q, \mathbb{R}^k))$  where  $\sigma, \tau \in \{-0, 1, -1\}^{\mathcal{P}}$ , with  $\bar{C} \cap D \neq \emptyset$ , there exists  $(p, q, \gamma) \in N$ , such that  $q \in D$  and  $\gamma \setminus \{q\} \in C$ .

Thus, if  $C$  and  $D$  are two semi-algebraically connected components of two distinct sign conditions whose union is semi-algebraically connected then there exists  $(p, q, \gamma) \in N$  such that  $\gamma$  connects the point  $p \in C$  with the point  $q \in D$  through a semi-algebraic path lying in  $C \cup D$ .

We first describe the algorithm constructing data for adjacencies and then prove its correctness.

### *Algorithm 16.7.* [Data for Adjacencies]

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ , and  $\text{Zer}(Q, \mathbb{R}^k)$  is of real dimension  $k'$ , a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** a set  $N$  of data for adjacencies for  $\mathcal{P}$  on  $\text{Zer}(Q, \mathbb{R}^k)$ , described by real univariate representations and parametrized real univariate representations.

- **Complexity:**  $s^{k'+1} d^{O(k)}$  where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**
  - Introduce a new variable  $\beta$  and define  $\mathcal{P}' = \{\{P, P + \beta, P - \beta\}, P \in \mathcal{P}\}$ .
  - Call Algorithm 13.3 (Sampling on an Algebraic Set) with input  $Q$  and  $\mathcal{P}'$  and structure  $D[\beta] \subset \mathbb{R}\langle\beta\rangle$  to obtain a set of real univariate representations.
  - For each associated point  $p(\beta)$ , compute  $q = \lim_{\beta} (p(\beta))$ , using Algorithm 12.14 (Limit of Bounded Points). The point  $p(\beta)$  is represented as a real  $k$ -univariate representation  $(u, \sigma)$  with

$$u = (f(\beta, T), g_0(\beta, T), \dots, g_k(\beta, T)).$$

Replacing  $\beta$  in  $u$  by a small enough  $t_0 \in \mathbb{R}$  using Algorithm 11.20 (Removal of Infinitesimals). Call Algorithm 11.20 (Removal of Infinitesimals) with input the polynomial  $f$  as well as the family of polynomials  $\{P_u | P \in \mathcal{P}\}$  (see Notation 13.8) to obtain  $t_0 \in \mathbb{R}$  replacing  $\beta$ . Letting  $t$  vary over the interval  $[0, t_0]$  gives a semi-algebraic path  $\gamma$  joining  $p(t_0)$  to  $q$ . Include the triple  $(q, p(t_0), \gamma)$  in the set  $N$ .

The proof of correctness uses the following lemma.

**Lemma 16.12.** *Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$  with  $\text{Zer}(Q, \mathbb{R}^k) \subset B(0, 1/c)$ . Let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a finite set of polynomials and*

$$\mathcal{P}' = \{\{P, P + \beta, P - \beta\}, P \in \mathcal{P}\}.$$

*Suppose that  $\sigma$  and  $\tau$  are distinct realizable sign conditions on  $\mathcal{P}$  and that  $C$  and  $D$  are two semi-algebraically connected components of  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$ , and  $\text{Reali}(\tau, \text{Zer}(Q, \mathbb{R}^k))$  respectively such that  $\bar{C} \cap D \neq \emptyset$ . Then there is a semi-algebraically connected component  $C'$ , of a realizable sign condition  $\sigma'$  of  $\mathcal{P}'$  on  $\text{Zer}(Q, \mathbb{R}^k)$  such that  $C' \subset \text{Ext}(C, \mathbb{R}\langle\beta\rangle)$  and  $\lim_{\beta}(C') \subset D$ .*

**Proof:** Let  $\mathcal{P} = \{P_1, \dots, P_s\}$ . Suppose without loss of generality that  $\sigma$  is

$$P_1 = \dots = P_{\ell} = 0, P_{\ell+1} > 0, \dots, P_s > 0.$$

After a possible re-ordering of the indices,  $\tau$  is

$$P_1 = \dots = P_m = 0, P_{m+1} > 0, \dots, P_s > 0$$

with  $m > \ell$ . This is clear since a point  $p \in \bar{C} \cap D$  must satisfy

$$P_1 = \dots = P_{\ell} = 0, P_{\ell+1} \geq 0, \dots, P_s \geq 0.$$

Consider the set defined by the formula  $\sigma'$

$$\begin{aligned} P_1 = \dots = P_{\ell} = 0, \\ 0 \leq P_{\ell+1} \leq \beta, \dots, 0 \leq P_m \leq \beta, \\ P_{m+1} > 0, \dots, P_s > 0. \end{aligned}$$

Let us prove first that the realization of  $\sigma'$  is non-empty. Let  $x \in \bar{C} \cap D$ . According to Theorem 3.19 (Curve Selection Lemma), there is a semi-algebraic path  $\gamma$  such that  $\gamma(0) = x$ ,  $\gamma((0, 1]) \subset C$ . Since at  $\gamma(1)$  we have  $P_{\ell+1} > 0, \dots, P_m > 0$  and at  $\gamma(0)$  we have  $P_{\ell+1} = \dots = P_m = 0$ , there exists  $t \in \mathbb{R} \langle \beta \rangle$  such that  $0 < P_{\ell+1} \leq \beta, \dots, 0 < P_m \leq \beta$  on  $\gamma((0, t])$  (use Exercise 3.1 part 3).

It is clear that the realization of  $\sigma'$  is contained in the extension of  $\sigma$  to  $\mathbb{R} \langle \beta \rangle$ . Consider the semi-algebraically connected component  $C'$  of the realization of  $\sigma'$  that contains  $y = \gamma(1)$ . It is clear that  $C' \subset \text{Ext}(C, \mathbb{R} \langle \beta \rangle)$ . Moreover,  $\lim_{\beta}(C')$  satisfies sign condition  $\tau$  and contains  $x \in D$ .

Since,  $\lim_{\beta}$  maps semi-algebraically connected sets to semi-algebraically connected sets, we see that  $\lim_{\beta}(C) \subset D$ . □

**Proof of correctness:** We need to show that the set of triples computed above is set of data for adjacencies for  $\mathcal{P}$ . This is an immediate consequence of Lemma 16.12.

It is clear that if  $p(\beta)$  is a point in  $C'$   $q = \lim_{\beta}(p)$ , and  $\gamma$  is the semi-algebraic path obtained by replacing  $\beta$  by a small enough  $t > t_0$  in  $p(\varepsilon)$  then  $p(t_0) \in C$ ,  $q \in D$  and  $\gamma$  is a semi-algebraic path joining  $p(t_0)$  and  $q$  contained in  $C$  except at the endpoint  $q$ . □

**Complexity analysis:** The complexity of the whole computation is  $s \sum_{j \leq k'} 4^j \binom{s}{j} d^{O(k)}$  in  $D[\beta]$ , using the complexity analyses of Algorithm 13.3 (Sampling on a Bounded Algebraic Set) (with the extra remark that  $P, P - \beta$  and  $P + \beta$  have no common zeroes) and Algorithm 11.20 (Removal of Infinitesimals). Since the degree in  $\beta$  of the intermediate computations is also bounded by  $d^{O(k)}$ , the complexity in  $D$  is finally  $s \sum_{j \leq k'} 4^j \binom{s}{j} d^{O(k)} = s^{k'+1} d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k)}$ . □

We can now describe the semi-algebraically connected components of a semi-algebraic set.

*Algorithm 16.8. [Connected Components of a Semi-algebraic Set]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:** a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ , a  $\mathcal{P}$ -semi-algebraic set  $S_a$ .
- **Output:** a description of the semi-algebraically connected components of  $S$ .
- **Complexity:**  $s^{k+1} d^{O(k^4)}$  where  $s$  is a bound on the number of polynomials in  $\mathcal{P}$  and  $d$  is a bound on their degree.

- **Procedure:** Using Algorithms 16.7 (Data for Adjacencies) and 16.6 (Basic Connected Components), compute the equivalence classes of the transitive closure of the adjacency relation between semi-algebraically connected components of the realizations of realizable sign condition, and take the union of the corresponding equivalence classes.

**Proof of correctness:** Follows from the correctness of Algorithms 16.7 (Data for Adjacencies) and 16.6 (Basic Connected Components).  $\square$

**Complexity analysis:** The complexity of the algorithm is bounded by  $s^{k+1}d^{O(k^4)}$  using the preceding results on the complexity of Algorithm 16.7 (Data for Adjacencies) and Algorithm 16.6 (Basic Connected Components). The degrees of the polynomials that appear in the output are bounded by  $d^{O(k^3)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^4)}$ .  $\square$

We have proved the following theorem:

**Theorem 16.13.** *Let  $\mathcal{P} = \{P_1, \dots, P_s\} \subset D[X_1, \dots, X_k]$  with  $\deg(P_i) \leq d$ , for  $1 \leq i \leq s$  and a semi-algebraic set  $S$  defined by a  $\mathcal{P}$  quantifier-free formula. There exists an algorithm that outputs quantifier-free semi-algebraic descriptions of all the semi-algebraically connected components of  $S$ . The complexity of the algorithm is bounded by  $s^{k+1}d^{O(k^4)}$ . The degrees of the polynomials that appear in the output are bounded by  $d^{O(k^3)}$ . Moreover, if the input polynomials have integer coefficients whose bitsize is bounded by  $\tau$  the bitsize of coefficients output is  $\tau d^{O(k^3)}$ .*

## 16.5 Roadmap Algorithm

Our aim in this section is to construct a roadmap of a semi-algebraic defined by a  $\mathcal{P}$ -quantifier free formula on an algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  of dimension  $k'$ . We use the construction of approximating varieties described in Section 13.3 in order to achieve better complexity for our algorithm.

Let  $S$  be an arbitrary semi-algebraic set defined by a finite set of polynomials  $\mathcal{P}$  which is contained in a bounded algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  of real dimension  $k'$ .

We first assume that  $\text{Zer}(Q, \mathbb{R}^k)$  is bounded. The idea is to construct uniform roadmaps for a perturbed finite set of polynomials which are in general position over approximating varieties (see Chapter 13 page 523) which are close to  $\text{Zer}(Q, \mathbb{R}^k)$  and of dimension  $k'$ .

We then take the limits of the curves obtained when the parameter of deformation tends to 0, i.e.the images of the curves so constructed under a  $\lim$  map.

We first describe this limit process. The idea is to modify Algorithm 15.2 (Curve Segments) so that the limit of the curve segments when the parameter of deformation tends to 0 is also output.

*Algorithm 16.9. [Modified Curve Segments]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, X_2, \dots, X_k]$ , such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$ ,
  - $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$
  - a polynomial  $\bar{Q} \in D[\varepsilon, X_1, X_2, \dots, X_k]$ , for which

$$\lim_{\varepsilon} (\text{Zer}(\bar{Q}, R\langle\varepsilon\rangle^k) \subset \text{Zer}(Q, R^k)$$

- a triangular Thom encoding  $\mathcal{T}, \sigma$  specifying  $z \in R\langle\varepsilon\rangle^{i-1}$  with coefficients in  $D[\varepsilon]$ ,
- a triangular Thom encoding  $\mathcal{T}', \sigma'$  specifying  $\lim_{\varepsilon}(z) \in R^{i-1}$ , with coefficients in  $D$
- a set of at most  $m$  points,  $\mathcal{N} \subset \text{Zer}(\bar{Q}, R\langle\varepsilon\rangle^k)$ , where each point of  $\mathcal{N}$  is defined by a real  $k$ -univariate representation  $u, \sigma$  with coefficients in  $D[\varepsilon]$ , above  $\mathcal{T}, \sigma$ .

**Output:**

- An ordered list of Thom encodings  $A_1, \alpha_1, \dots, A_{\ell}, \alpha_{\ell}$  above  $\mathcal{T}, \sigma$  specifying points  $(z, v_1), \dots, (z, v_{\ell})$  with  $v_1 < \dots < v_{\ell}$ .
- An ordered list of Thom encodings  $B_1, \beta_1, \dots, B_{\ell}, \beta_{\ell}$  above  $\mathcal{T}', \sigma'$  specifying the image under  $\lim_{\varepsilon}$  of these distinguished values:
- For every  $j = 1, \dots, \ell$ ,
  - a finite set  $\mathcal{D}_j$  of real univariate representation above  $\mathcal{T}, A_j, \sigma, \alpha_j$ . The associated points are called distinguished points.
  - a finite set  $\mathcal{D}'_j$  of real univariate representation above  $\mathcal{T}', B_j, \sigma', \beta_j$ . The associated points are the image under  $\lim_{\varepsilon}$  of the distinguished points of  $\mathcal{D}_j$ .
  - a finite set  $\mathcal{C}_j$  of curve segment representations above

$$\mathcal{T}, \sigma, A_j, \alpha_j, A_{j+1}, \alpha_{j+1}.$$

The associated curve segments are called distinguished curves.

- a finite set  $\mathcal{C}'_j$  of curve segment representations

$$\mathcal{T}', \sigma', B_j, \beta_j, B_{j+1}, \beta_{j+1}.$$

with associated curve segments the image under  $\lim_{\varepsilon}$  of the curve segments in  $\mathcal{C}_j$ .

- a list of pairs of elements of  $\mathcal{C}_j$  and  $\mathcal{D}_j$  (resp.  $\mathcal{C}_{j+1}$  and  $\mathcal{D}_j$ ) describing the adjacency relations between distinguished curves and distinguished points.

The distinguished curves and points are contained in  $\text{Zer}(\bar{Q}, \mathbb{R}\langle\varepsilon\rangle^k)_z$ . Among the distinguished values are the first coordinates of the points in  $\mathcal{N}$  as well as the pseudo-critical values of  $\text{Zer}(\bar{Q}, \mathbb{R}\langle\varepsilon\rangle^k)_z$ . The sets of distinguished values, distinguished curves and distinguished points satisfy the following properties.

- CS<sub>1</sub>: For every  $v \in \mathbb{R}\langle\varepsilon\rangle$  the set of distinguished curve and distinguished points output intersect every semi-algebraically connected component of  $\text{Zer}(\bar{Q}, \mathbb{R}\langle\varepsilon\rangle^k)_{z,v}$ .
- CS<sub>2</sub>: For each distinguished curve output over an interval with endpoint a given distinguished value, there exists a distinguished point over this distinguished value which belongs to the closure of the curve segment.
- **Complexity:**  $d^{O(ik)}$ , where  $d$  is a bound on the degree of  $Q$  and  $\bar{Q}$ , and  $O(d)^k$  is a bound on the degree of the polynomials in  $\mathcal{T}$ , on the degree of the univariate representations in  $\mathcal{N}$  and on the number of these univariate representations.

- **Procedure:**

- Step 1: Perform Algorithm 12.10 (Parametrized Multiplication Table) with input  $\overline{\text{Cr}}(\bar{Q}^2, \zeta)$ , (using Notation 12.46) and parameter  $X_{\leq i}$ . Perform Algorithm 12.15 (Parametrized Limit of Bounded Points) and output  $\mathcal{U}$ .
- Step 2: For every  $u, \tau \in \mathcal{N}$ , compute  $\text{proj}_i(u), \text{proj}_i(\tau)$  using Algorithm 15.1 (Projection), add to  $\mathcal{D}$  the polynomial  $\text{proj}_i(u)$ .
- Step 3: Compute the Thom encodings of the zeroes of  $\mathcal{T}, A, A \in \mathcal{D}$  above  $\mathcal{T}, \sigma$  using Algorithms 12.20 (Triangular Thom Encoding) and output their ordered list  $A_1, \alpha_1, \dots, A_\ell, \alpha_\ell$  and the corresponding ordered list  $v_1 < \dots < v_\ell$  of distinguished values using 12.21.

Compute the Thom encoding of  $\lim_\varepsilon(v_1) \leq \dots \leq \lim_\varepsilon(v_\ell)$ .

- Step 4: For every  $j = 1, \dots, \ell$  and every  $(f, g_0, g_i, \dots, g_k), \tau \in \mathcal{N}$  such that  $\text{proj}_i(\tau) = \alpha_j$ , append  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{D}_j$ .
- Step 5: For every  $j = 1, \dots, \ell$  output a finite set of univariate representations  $\mathcal{D}_j$  such that the set of associated points contains the set of  $X_i$ -pseudo-critical points of  $\text{Zer}(\bar{Q}, \mathbb{R}\langle\varepsilon\rangle^k)_{v_i}$  as well as a set of univariate representations  $\mathcal{D}'_j$  with associated points the  $\lim_\varepsilon$  image of the points associated to  $\mathcal{D}_j$ .

For every  $j = 1, \dots, \ell$  and every  $u = (f, g_0, g_i, \dots, g_k) \in \mathcal{U}$ , compute the Thom encodings  $\tau$  of the roots of  $\mathcal{T}, f$  such that  $\text{proj}_i(\tau) = \alpha_j$ , using Algorithm 12.20 (Triangular Thom Encoding) and append all pairs  $(f, g_0, g_{i+1}, \dots, g_k), \tau$  to  $\mathcal{D}_j$  when the corresponding associated point belongs to  $\text{Zer}(\bar{Q}, \mathbb{R}\langle\varepsilon\rangle^k)_z$ .

For every  $u \in \mathcal{D}_j$ , such that  $o(f) = o(u)$ , put

$$\hat{u}(X_{\leq i}, T) = \lim_{\varepsilon} (\varepsilon^{-o(f)} u(\varepsilon, X_{\leq i}, T)).$$

with coefficients in  $\mathbb{D}[X_{\leq i}, T]$  in  $\mathcal{D}'_j$ .

- Step 6: Output on each open interval  $(v_j, v_{j+1})$  a finite set of curve segments  $\mathcal{C}_j$  such that for every  $v \in (v_j, v_{j+1})$  the set of associated points contains the set of  $X_i$ -pseudo-critical points of  $\text{Zer}(Q, \mathbb{R}^k)_v$ .

For every  $j = 1, \dots, \ell - 1$  and every  $u = (f, g_0, g_{i+1}, \dots, g_k) \in \mathcal{U}$ , compute the Thom encodings  $\rho$  of the roots of  $f(z, v, T)$  over  $(v_j, v_{j+1})$  using Algorithm 12.22 (Triangular Intermediate Points) and Algorithm 12.20 (Triangular Thom Encoding). Append pairs  $u, \rho$  to  $\mathcal{C}_j$  when the corresponding associated curve is included in  $\text{Zer}(\bar{Q}, \mathbb{R}\langle \varepsilon \rangle^k)_z$ .

For every  $u \in \mathcal{C}_j$ , such that  $o(f) = o(u)$ , put

$$\bar{u}(X_{\leq i}, T) = \lim_{\varepsilon} (\varepsilon^{-o(f)} u(\varepsilon, X_{\leq i}, T))$$

with coefficients in  $\mathbb{D}[X_{\leq i}, T]$  in  $\mathcal{C}'_j$ .

- Step 7: Determine adjacencies between curve segments and points. For every point of  $\mathcal{D}_j$  specified by

$$v' = (p, q_0, q_{i+1}, \dots, q_k), \tau', \{p, q_0, q_{i+1}, \dots, q_k\} \subset \mathbb{D}[X_{\leq i}][T]$$

and every curve segment representation of  $\mathcal{C}_j$  specified by

$$v = (f, g_0, g_{i+1}, \dots, g_k), \tau, \{f, g_0, g_{i+1}, \dots, g_k\} \subset \mathbb{D}[X_{\leq i}][T],$$

decide whether the corresponding point  $t$  is adjacent to the corresponding curve segment as follows: compute the first  $\nu$  such that  $\partial^\nu g_0 / \partial X_i^\nu(v_j, t)$  is not zero and decide whether for every  $\ell = i + 1, \dots, k$

$$\frac{\partial^\nu g_\ell}{\partial X_i^\nu}(v_j, t) q_0(t) - \frac{\partial^\nu g_0}{\partial X_i^\nu}(v_j, t) q_\ell(t)$$

is zero. This is done using Algorithm 12.19 (Triangular Sign Determination).

Repeat the same process for every element of  $\mathcal{D}_{j+1}$  and every curve segment representation of  $\mathcal{C}_j$ .

**Proof of correctness:** It follows from Proposition 12.42, the correctness of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination), Algorithm 15.1 (Projection), Algorithm 12.22 (Triangular Intermediate Points), Algorithm 12.20 (Triangular Thom Encoding), Algorithm 12.21 (Triangular Comparison of Roots) and Algorithm 12.19 (Triangular Sign Determination). □

**Complexity analysis:** Step 1: This step requires  $d^{O(i(k-i))}$  arithmetic operations in  $\mathbb{D}$ , using the complexity analysis of Algorithm 12.10 (Parametrized Multiplication Table), Algorithm 12.15 (Parametrized Limit of Bounded Points), Algorithm 11.19 (Restricted Elimination). There are  $d^{O(k-i)}$  parametrized univariate representations computed in this step and each polynomial in these representations has degree  $O(d)^{k-i}$ .

Step 2: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 15.1 (Projection).

Step 3: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 12.20 (Triangular Thom Encoding).

Step 4: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

Step 5: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 12.20 (Triangular Thom Encoding).

Step 6: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 12.22 (Triangular Intermediate Points), Algorithm 12.20 (Triangular Thom Encoding).

Step 7: This step requires  $d^{O(ik)}$  arithmetic operations in  $D$ , using the complexity analysis of Algorithm 12.19 (Triangular Sign Determination).

Thus, the complexity is  $d^{O(ik)}$ . The number of distinguished values is bounded by  $d^{O(k)}$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(ik)}$ .  $\square$

We describe the construction of a set,  $L$ , such  $L$  meets every semi-algebraically connected component of every realizable weak sign condition of  $\mathcal{P}$  on  $\lim_{\eta}(Z_j) \cap \lim_{\eta}(Z_{\ell})$ , where  $Z_j$  and  $Z_{\ell}$  are the approximating varieties defined in Notation 13.30, for every  $0 \leq j \leq k'(k - k')$  and  $0 \leq \ell \leq k'(k - k')$ ,

*Algorithm 16.10.* [Linking Points]

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$  is of real dimension  $k'$ ,
  - a finite set  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** a set of points  $L$  such that for every  $0 \leq j \leq k'(k - k')$  and  $0 \leq \ell \leq k'(k - k')$ ,  $L$  meets every semi-algebraically connected component of every realizable weak sign condition of  $\mathcal{P}$  on  $\lim_{\eta}(Z_j) \cap \lim_{\eta}(Z_{\ell})$ .
- **Complexity:**  $s^{k'+1} d^{O(k^2)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$ .
- **Procedure:**
  - For every  $0 \leq j \leq k'(k - k')$ , denote by  $\mathcal{R}_j$  the set of polynomials in  $k + 1$  variables obtained after two steps of Algorithm 14.1 (Block Elimination) applied to the polynomials appearing in the formula

$$(\exists (X, T)) \| (X, T) - Y \|^2 < Z^2 \wedge T > 0 \wedge \bigwedge_{P \in \text{App}(Q_j, T)} P(X) = 0$$



describing the closure of the set

$$\{(x, t) \in \mathbb{R}^{k+1} \mid t > 0 \wedge \bigwedge_{P \in \text{App}(Q_j, t)} P(x) = 0\},$$

in order to eliminate  $Z$  and  $X, T$ . Denote by  $\mathcal{P}_j$  the set of polynomials in  $k$  variables obtained by substituting 0 for  $T$  in  $M_{k',j}(\mathcal{R}_j)$  (see Notation 13.26).

- For every  $0 \leq j \leq k'(k - k')$  and  $0 \leq \ell \leq k'(k - k')$ , apply Algorithm 13.3 (Sampling on a Bounded Algebraic Set), with input  $\text{Zer}(Q, \mathbb{R}^k)$ ,  $\mathcal{P} \cup \mathcal{P}_\ell \cup \mathcal{P}_j$  to obtain the set  $L_{\ell,j}$ . The set  $L$  is the union of the  $L_{\ell,j}$ .

**Proof of correctness:** Note that,

$$\lim_{\eta} (\text{Zer}(\text{App}(Q_\ell, \eta), \mathbb{R}\langle \eta \rangle^k))$$

is the closure of

$$\{(x, t) \in \mathbb{R}^{k+1} \mid t > 0 \wedge \bigwedge_{P \in \text{App}(Q_\ell, t)} P(x) = 0\} \cap \{t = 0\}.$$

The polynomials  $\mathcal{R}_j$  have the property that the closure of

$$\{(x, t) \in \mathbb{R}^{k+1} \mid t > 0 \wedge \bigwedge_{P \in P \in \text{App}(Q_j, t)} P(x) = 0\}$$

is the union of semi-algebraically connected components of sets defined by sign conditions over  $\mathcal{R}_j$  (see page 556).

For every  $0 \leq j \leq k'(k - k')$  and  $1 \leq \ell \leq k'(k - k')$ ,  $L$  meets every semi-algebraically connected component of every realizable weak sign condition of  $\mathcal{P}$  on  $\lim_{\eta} (Z_j) \cap \lim_{\eta} (Z_\ell)$ . □

**Complexity analysis:** According to the complexity of Algorithm 14.1 (Block Elimination), the set  $\mathcal{R}_j$  has  $d^{O(k)}$  polynomials and each of these polynomials has degree at most  $d^{O(k)}$ .

According to the complexity of Algorithm 13.3 (Sampling on a Bounded Algebraic Set), the set  $L_{i,j}$  consists of  $\sum_{j=0}^{k'} \binom{s}{j} 4^j d^{O(k^2)}$  points defined by polynomials of degree at most  $d^{O(k^2)}$ . The complexity is

$$s \sum_{j=0}^{k'} \binom{s}{j} 4^j d^{O(k^2)} = s^{k'+1} d^{O(k^2)}.$$

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

In order to ensure that the roadmaps constructed on the various approximating varieties take into account connectivity in the original algebraic set, we need to add points in the various roadmaps for approximating varieties.

*Algorithm 16.11.* **[Touching Points]**

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$  is of real dimension  $k'$ ,
  - a real univariate representation  $u$  describing a point  $p \in \text{Zer}(Q, R^k)$ .
- **Output:** for every  $0 \leq j \leq k'(k - k')$  such that  $Z_j$  is infinitesimally close to  $p$ , a set of real univariate representations describing points meeting every semi-algebraically connected component of  $Z_j$  infinitesimally close to  $p$ .
- **Complexity:**  $s^{k'+1}d^{O(k^2)}$  where  $s$  is a bound on the number of elements of  $\mathcal{P}$ ,  $d$  is a bound on the degrees of  $Q$  and the elements of  $\mathcal{P}$  and  $d^{O(k^2)}$  is a bound on the degrees of  $u$ .
- **Procedure:**
  - Let  $u = (f, g_0, \dots, g_k), \sigma$ . For every  $0 \leq j \leq k'(k - k')$  proceed as follows. Let  $\beta$  be a new variable and let  $P_p(T, X_1, \dots, X_k)$  be the system

$$\{f(T), \sum_{i=1}^k (g_0(T)X_i - g_i(T))^2 - g_0(T)^2\beta^2\}$$

Call Algorithm 13.3 (Sampling on a Bounded Algebraic Set) with input  $\text{App}(Q_j, \eta)$ ,  $P_p$  and  $\text{Der}(f)$  in the ring  $D[\beta, \eta]$ . For each real univariate representation obtained, keep all those corresponding to points  $q$  at which the sign of  $P_p$  is negative and such that the sign condition satisfied by  $\text{Der}(f)$  at  $\lim_{\beta, \eta}(q)$  is  $\sigma$  and discard the rest. Denote by  $\mathcal{U}_j$  the real univariate representations representing points of  $Z_j$  obtained by applying  $M_{k', j}$  (see Notation 13.26) to the real univariate representation associated to  $q$ .

- Output the set  $\mathcal{U} = \bigcup_{j=0}^{k'(k-k')} \mathcal{U}_j$  of real univariate representations so obtained. The touching points are the points associated to the elements of  $\mathcal{U}$ .

**Proof of correctness:** Immediate. □

**Complexity analysis:** The number of arithmetic operations in  $D$  for computing the set of touching points is  $s \sum_{j=0}^{k'} \binom{s}{j} 4^j d^{O(k^2)} = s^{k'+1} d^{O(k^2)}$ . This follows from the complexity of Algorithm 13.3 (Sampling on a Bounded Algebraic Set).

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ . □

We now describe the roadmap algorithm in the bounded case.

*Algorithm 16.12. [Bounded Roadmap]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, R^k) \subset B(0, 1/c)$  is of real dimension  $k'$ ,
  - a semi-algebraic subset  $S$  of  $\text{Zer}(Q, R^k)$  defined by a  $\mathcal{P}$ -quantifier-free formula where  $\mathcal{P} \subset D[X_1, \dots, X_k]$ .
- **Output:** a roadmap for  $S$ .
- **Complexity:**  $s^{k'+1}d^{O(k^2)}$  where  $s$  is a bound on the number of elements of  $\mathcal{P}$ , and  $d$  is a bound on the degree of  $Q$  and of the polynomials in  $\mathcal{P}$ .
- **Procedure:**
  - Let  $d' = 2(d+1)$ .
  - For every  $0 \leq \ell \leq k'(k-k')$ , define

$$\bar{Q}_\ell = Q_\ell^2 + (\varepsilon^2 (X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 1)^2,$$

and define  $\text{App}(\bar{Q}_\ell, \eta)$  and  $\mathcal{P}_\ell^*$ , using Notation 13.30 and Notation 13.32. Use a modified version of Algorithm 16.4 (Uniform Roadmap) with input  $(\text{App}(\bar{Q}_\ell, \eta), \mathcal{P}_\ell^*)$  using Algorithm 16.9 (Modified Curve Segments) rather than Algorithm 15.2 (Curve Segments).

- Call Algorithm 16.7 (Data for Adjacencies) and Algorithm 16.10 (Linking Points). For each element of  $\bar{N} \cup L$ , obtained above apply Algorithm 16.11 (Touching points). This defines a set  $\mathcal{A}_\ell$  of real univariate representations. Connect the points associated to the elements of  $\mathcal{A}_\ell$  to the uniform roadmap for  $(\text{App}(\bar{Q}_\ell, \eta), \mathcal{P}_\ell^*)$  using a modified version of Algorithm 16.2 (Bounded Connecting), using Algorithm 16.9 (Modified Curve Segments) rather than Algorithm 15.2 (Curve Segments).
- Output the image of the segments and points constructed above under the  $\lim_{\gamma, \eta}$  map, using the computation done in the calls to Algorithm 16.9 (Modified Curve Segments) and retain only those portions which are in the given set  $S$ .

**Proof of correctness:** The correctness follows from the correctness of Algorithm 16.4 (Uniform Roadmap), Algorithm 16.9 (Modified Curve Segments), Algorithm 16.7 (Data for Adjacencies), Algorithm 16.10 (Linking Points) and Algorithm 16.11 (Touching points), as well as Proposition 13.33 and Proposition 13.35.  $\square$

**Complexity analysis:** The number of arithmetic operations for computing the set of added points  $\mathcal{A}_\ell$  is  $s \sum_{j=0}^{k'} \binom{s}{j} 4^j d^{O(k^2)}$  in  $D$ , using the complexity analysis of Algorithm 16.7 (Data for Adjacencies), Algorithm 16.10 (Linking Points) and Algorithm 16.11 (Touching points).

Since the set  $\mathcal{P}_\ell^*$  is in  $k'$ -general position with respect to  $\text{App}(\bar{Q}_\ell, d', \varepsilon, \eta)$  according to Proposition 13.33, using the complexity bound of Algorithm 16.4 (Uniform Roadmap), we see that the complexity is bounded by  $s^{k'+1}d^{O(k^2)}$  in  $D$ .

Similarly, using the complexity bounds for Algorithm 16.2 (Bounded Connecting), the complexity of connecting a point  $x$  described by polynomials of degree at most  $d^{O(k)}$  to the roadmap is  $k' s d^{O(k^2)}$  in  $D$ .

If  $D = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .  $\square$

Now we show how to modify Algorithm 16.12 (Bounded Roadmap) to handle the case when the input algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  is not bounded.

*Algorithm 16.13. [General Roadmap]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .
- **Input:**
  - a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k)$  is of real dimension  $k'$ ,
  - a semi-algebraic subset  $S$  of  $\text{Zer}(Q, \mathbb{R}^k)$  described by a finite set  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ .
- **Output:** a roadmap for  $S$ .
- **Complexity:**  $s^{k'+1}d^{O(k^2)}$  where  $s$  is a bound on the number of elements of  $\mathcal{P}$ , and  $d$  is a bound on the degree of  $Q$  and of the polynomials in  $\mathcal{P}$ .
- **Procedure:**
  - Step 1: Introduce new variables  $X_{k+1}$  and  $\varepsilon$  and replace  $Q$  by the polynomial  $Q^* = Q^2 + (\varepsilon^2 (X_1^2 + \dots + X_{k+1}^2) - 1)^2$ . Let  $S^* \in \mathbb{R}\langle \varepsilon \rangle^{k+1}$  be the set defined by the same formula as  $S$  but with  $Q$  replaced by  $Q^*$ . Run Algorithm 16.12 (Bounded Roadmap) with input  $Q^*$  and  $S^*$  and output a roadmap for  $\text{RM}(S^*)$ , composed of points and curves whose description involves  $\varepsilon$ .
  - Step 2: Denote by  $\mathcal{L}$  be the set of all polynomials in  $D[\varepsilon]$  whose signs were determined in the various calls to the Multivariate Sign Determination Algorithm in Step 1. Replace  $\varepsilon$  by

$$a = \min_{P \in \mathcal{L}} c(P)$$

(Definition 10.5) in the output roadmap to obtain a roadmap  $\text{RM}(S_a)$ . When projected on  $\mathbb{R}^k$ , this gives a roadmap  $\text{RM}(S) \cap B(0, 1/a)$ .

- Step 3: Collect all the points  $(y_1, \dots, y_k)$  in the roadmap which satisfies  $\varepsilon^2 (y_1^2 + \dots + y_k^2) = 1$ . Each such point is described by a univariate representation involving  $\varepsilon$ . Add to the roadmap the curve segment obtained by treating  $\varepsilon$  as a parameter and letting  $\varepsilon$  vary over  $(0, a, ]$ , to get a roadmap  $\text{RM}(S)$ .

**Proof of correctness:** Follows from the correctness of Algorithm 16.12 (Bounded Roadmap).  $\square$

**Complexity analysis:** The complexity is bounded by  $s^{k'+1}d^{O(k^2)}$  in  $\mathbb{D}$  and coincides with the complexity of Algorithm 16.12 (Bounded Roadmap).

If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .  $\square$

Using the preceding algorithms we can now prove.

**Theorem 16.14.** *Let  $Q \in \mathbb{R}[X_1, \dots, X_k]$  with  $\text{Zer}(Q, \mathbb{R}^k)$  of dimension  $k'$  and let  $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$  be a set of at most  $s$  polynomials for which the degrees of the polynomials in  $\mathcal{P}$  and  $Q$  are bounded by  $d$ . Let  $S$  be a semi-algebraic subset of  $\text{Zer}(Q, \mathbb{R}^k)$  defined by a  $\mathcal{P}$ -quantifier-free formula.*

- a) *Let  $p \in \text{Zer}(Q, \mathbb{R}^k)$  a point which is represented by a  $k$ -univariate representation with specified Thom encoding  $(u, \sigma)$  of degree  $d^{O(k)}$ . There is an algorithm whose output is a semi-algebraic path connecting  $p$  to  $\text{RM}(S)$ . The complexity of the algorithm in the ring  $\mathbb{D}$  generated by the coefficients of  $Q$ ,  $u$  and the elements of  $\mathcal{P}$  is bounded by  $k's d^{O(k^2)}$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .*
- b) *There is an algorithm whose output is exactly one point in every semi-algebraically connected component of  $S$ . The complexity in the ring generated by the coefficients of  $Q$  and  $\mathcal{P}$  is bounded by  $s^{k'+1}d^{O(k^2)}$ . In particular, this algorithm counts the number semi-algebraically connected component of  $S$  in time  $s^{k'+1}d^{O(k^2)}$  in the ring  $\mathbb{D}$  generated by the coefficients of  $Q$  and the coefficients of the elements of  $\mathcal{P}$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .*
- c) *Let  $p$  and  $q$  be two points that are represented by real  $k$ -univariate real representation  $u$  and  $v$ , of degree  $d^{O(k)}$  belonging to  $S$ . There is an algorithm deciding whether  $p$  and  $q$  belong to the same connected component of  $S$ . The complexity in the ring  $\mathbb{D}$  generated by the coefficients of  $Q$ ,  $u$ ,  $v$  and the coefficients of the polynomials in  $\mathcal{P}$ . is bounded by  $s^{k'+1}d^{O(k^2)}$ . If  $\mathbb{D} = \mathbb{Z}$ , and the bitsizes of the coefficients of the polynomials are bounded by  $\tau$ , then the bitsizes of the integers appearing in the intermediate computations and the output are bounded by  $\tau d^{O(k^2)}$ .*

**Proof:** a) In order to connect a point  $x$  to the roadmap in the bounded case, chose a  $0 \leq j \leq k'(k - k')$  such that  $x \in \lim_{\gamma, \eta} (Z_j)$  and construct a point  $x_j$  infinitesimally close to  $x$  in  $Z_j$  using Algorithm 13.3 (Sampling on an Algebraic Set) and Algorithm 16.11 (Touching Points). This point  $x_j$  is connected to the uniform roadmap  $\text{RM}(\text{App}(Q_\ell, \eta), \mathcal{P}_j^*)$  using a modified version of Algorithm 16.2 (Bounded Connecting) using Algorithm 16.9 (Modified Curve Segments) instead of Algorithm 15.2 (Curve Segments). Then output the image of the connecting curves under the map  $\lim_{\eta, \gamma}$  using the computations done in the calls to Algorithm 16.9 (Modified Curve Segments). In the unbounded case, we modify the preceding method using the same method as in Step 3 of Algorithm 16.13 (General Roadmap).

b) and c) are clear after a). □

## 16.6 Computing the First Betti Number of Semi-algebraic Sets

Our aim in this section is to compute the first Betti number of a  $\mathcal{P}$ -closed semi-algebraic set.

We first describe an algorithm for computing closed a contractible coverings of a  $\mathcal{P}$ -closed semi-algebraic set in single exponential time when the family  $\mathcal{P}$  is in general position. This algorithm computes parametrized connecting paths using Algorithm 16.5 (Parametrized Bounded) and uses them to construct a contractible covering.

We are given a polynomial  $Q \in D[X_1, \dots, X_k]$  such that  $\text{Zer}(Q, \mathbb{R}^k)$  is bounded and a finite set of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$  in strong  $k$ -general position with respect to  $Q$ .

We fix a closed semi-algebraic set  $S$  contained in  $\text{Zer}(Q, \mathbb{R}^k)$ . We follow the notations of Algorithm 16.5. and let  $\#\mathcal{A} = t$ . We denote by  $\text{SIGN}(S)$  the set of realizable sign conditions of  $\mathcal{A}$  on  $\text{Zer}(Q, \mathbb{R}^k)$  whose realizations are contained in  $S$ , remembering that  $\mathcal{P} \subset \mathcal{A}$ . For each  $\sigma \in \text{SIGN}(S)$   $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$  is contained in  $\text{Reali}(\theta, \text{Zer}(Q, \mathbb{R}^k))$  for some  $\theta \in \Theta$ . We denote by  $\gamma(\sigma)$  the restriction of  $\gamma(\theta)$  to the base  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$ . Since  $\gamma(\theta)$  is a parametrized path,  $\gamma(\sigma)$  is also a parametrized path. However, since  $\text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}^k))$  is not necessarily closed and bounded, we cannot use Proposition 15.15, and  $\text{Im}\gamma(\sigma)$  might not be contractible. In order to ensure contractibility, we restrict the base of  $\gamma(\sigma)$  to a slightly smaller set which is closed, using infinitesimals.

We introduce infinitesimals

$$\varepsilon_{2t} \gg \varepsilon_{2t-1} \gg \dots \gg \varepsilon_2 \gg \varepsilon_1 > 0.$$

For  $i = 1, \dots, 2t$  we denote by  $D_i$  the ring  $D[\varepsilon_{2t}, \dots, \varepsilon_i]$ , and by  $R_i$  the field  $R\langle \varepsilon_{2t} \rangle \dots \langle \varepsilon_i \rangle$ .

For  $\sigma \in \text{SIGN}(S)$  we define the level of  $\sigma$  by,

$$\text{level}(\sigma) = \#\{P \in \mathcal{A} \mid \sigma(P) = 0\}.$$

Given  $\sigma \in \text{SIGN}(S)$ , with  $\text{level}(\sigma) = j$ , we denote by  $\text{Reali}(\sigma_-)$  the set defined on  $\text{Zer}(Q, \mathbb{R}_{2j}^k)$  by the formula  $\sigma_-$  obtained by taking the conjunction of

$$\begin{aligned} P = 0, & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 0, \\ P \geq \varepsilon_{2j}, & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 1, \\ P \leq -\varepsilon_{2j}, & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = -1. \end{aligned}$$

Notice that  $\text{Reali}(\sigma_-) \subset \text{Reali}(\sigma, \text{Zer}(Q, \mathbb{R}_{2j}^k))$  is closed and bounded. Proposition 15.15 implies,

**Proposition 16.15.** *The set  $\gamma(\sigma)(\text{Reali}(\sigma_-))$  is semi-algebraically contractible.*

Note that the sets  $\gamma(\sigma)(\text{Reali}(\sigma_-))$  do not necessarily cover  $S$ . So we are going to enlarge them, preserving contractibility, to obtain a covering of  $S$ .

Given  $\sigma \in \text{SIGN}(S)$ , with  $\text{level}(\sigma) = j$ , we denote by  $\text{Reali}(\sigma^+)$  the set defined on  $\text{Zer}(Q, \mathbb{R}_{2j-1}^k)$ , by the formula  $\sigma^+$  obtained by taking the conjunction of

$$\begin{aligned} -\varepsilon_{2j-1} \leq P \leq \varepsilon_{2j-1} & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 0, \\ P \geq \varepsilon_{2j}, & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = 1, \\ P \leq -\varepsilon_{2j}, & \quad \text{for each } P \in \mathcal{A} \text{ such that } \sigma(P) = -1. \end{aligned}$$

with the formula  $\phi$  defining  $S$ . Let  $C(\sigma)$  be the set defined by,

$$C(\sigma) = \gamma(\sigma)(\text{Reali}(\sigma_-) \cup \text{Reali}(\sigma^+)).$$

We now prove that

**Proposition 16.16.**  *$C(\sigma)$  is semi-algebraically contractible.*

Let  $C$  be a closed and bounded semi-algebraic set contained in  $\mathbb{R}\langle\varepsilon\rangle^k$ . We can suppose without loss of generality that  $C$  is defined over  $\mathbb{R}[\varepsilon]$  by Proposition 2.82. We denote by  $C(t)$  the semi-algebraic subset of  $\mathbb{R}^k$  defined by replacing  $\varepsilon$  by  $t$  in the definition of  $C$ . Note that  $C(\varepsilon)$  is nothing but  $C$ .

We are going to use the following lemma.

**Lemma 16.17.** *Let  $B$  be a closed and bounded semi-algebraic set contained in  $\mathbb{R}^k$  and let  $C$  be a closed and bounded semi-algebraic set contained in  $\mathbb{R}\langle\varepsilon\rangle^k$ . If there exists  $t_0$  such that for every  $t < t' < t_0$ ,  $C(t) \subset C(t')$ , and  $\lim_\varepsilon(C) = B$ , then  $\text{Ext}(B, \mathbb{R}\langle\varepsilon\rangle)$  has the same homotopy type as  $C$ .*

**Proof:** Hardt's Triviality Theorem (Theorem 5.46) implies that there exists  $t_0 > 0$ , and a homeomorphism

$$\phi_{t_0}: C(t_0) \times (0, t_0] \rightarrow \cup_{0 < t \leq t_0} C t$$

which preserves  $C(t_0)$ . Replacing  $t_0$  by  $\varepsilon$  gives a homeomorphism

$$\phi(\varepsilon): C \times (0, \varepsilon] \rightarrow \cup_{0 < t \leq \varepsilon} C(t).$$

Defining

$$\psi: C \times [0, \varepsilon] \rightarrow C$$

by

$$\begin{aligned} \psi(x, s) &= \pi_{1\dots k} \circ \phi(x, s), & \text{if } s > 0 \\ \psi(x, 0) &= \lim_{s \rightarrow 0^+} \pi_{1\dots k} \circ \phi(x, s), \end{aligned}$$

it is clear that  $\psi$  is a semi-algebraic retraction of  $C$  to  $\text{Ext}(B, \mathbb{R}(\varepsilon))$ . □

We now prove Proposition 16.16.

**Proof of Proposition 16.16:** Apply Lemma 16.17 to  $C_\sigma$  and

$$\text{Ext}(\gamma(\sigma)(\text{Reali}(\sigma_-)), \mathbb{R}_{2j-1}):$$

thus  $C(\sigma)$  can be semi-algebraically retracted to  $\text{Ext}(\gamma(\sigma)(\text{Reali}(\sigma_-)), \mathbb{R}_{2j-1})$ .

Since  $\text{Ext}(\gamma(\sigma)(\text{Reali}(\sigma_-)), \mathbb{R}_{2j-1})$  is semi-algebraically contractible, so is  $C(\sigma)$ . □

We now prove that the sets  $\text{Ext}(C(\sigma), \mathbb{R}_1)$  form a covering of  $\text{Ext}(S, \mathbb{R}_1)$ .

**Proposition 16.18. [Covering property]**

$$\text{Ext}(S, \mathbb{R}_1) = \bigcup_{\sigma \in \text{SIGN}(S)} \text{Ext}(C(\sigma), \mathbb{R}_1).$$

The proposition is an immediate consequence of the following stronger result.

**Proposition 16.19.**

$$\text{Ext}(S, \mathbb{R}_1) = \bigcup_{\sigma \in \text{SIGN}(S)} \text{Reali}(\sigma_-^+, \mathbb{R}_1^k).$$

**Proof:** By definition,

$$\text{Ext}(S, \mathbb{R}_1) \supset \bigcup_{\sigma \in \text{SIGN}(S)} \text{Reali}(\sigma_-^+, \mathbb{R}_1^k).$$

We now prove the reverse inclusion. Clearly, we have that

$$S = \bigcup_{\sigma \in \text{SIGN}(S)} \text{Reali}(\sigma, \mathbb{R}).$$



Let  $x \in \text{Ext}(S, R_1)$  and  $\sigma$  be the sign condition of the family  $\mathcal{A}$  at  $x$  and let  $\text{level}(\sigma) = j$ . If  $x \in \text{Reali}(\sigma^+, R_1^k)$ , we are done. Otherwise, there exists  $P \in \mathcal{A}$ , such that  $x$  satisfies either  $0 < P(x) < \varepsilon_{2j}$  or  $-\varepsilon_{2j} < P(x) < 0$ . Let  $\mathcal{B} = \{P \in \mathcal{A} \mid \lim_{\varepsilon_{2j}} P(x) = 0\}$ . Clearly  $\#\mathcal{B} = j' > j$ . Let  $y = \lim_{\varepsilon_{2j}} x$ . Since,  $\text{Ext}(S, R_1)$  is closed and bounded and  $x \in \text{Ext}(S, R_1)$ ,  $y$  is also in  $\text{Ext}(S, R_1)$ . Let  $\tau$  be the sign condition of  $\mathcal{A}$  at  $y$  with  $\text{level}(\tau) = j' > j$ . If  $x \in \text{Reali}(\tau^+, R_1^k)$  we are done. Otherwise, for every  $P \in \mathcal{A}$  such that  $P(y) = 0$ , we have that  $-\varepsilon_{2j'-1} \leq P(x) \leq \varepsilon_{2j'-1}$ , since  $\lim_{\varepsilon_{2j}} (P(x)) = P(y) = 0$  and  $\varepsilon_{2j'-1} \gg \varepsilon_{2j}$ . So there exists  $P \in \mathcal{A}$  such that  $x$  satisfies either  $0 < P(x) < \varepsilon_{2j'}$  or  $-\varepsilon_{2j'} < P(x) < 0$ , and we replace  $\mathcal{B}$  by  $\{P \in \mathcal{A} \mid \lim_{\varepsilon_{2j'}} P(x) = 0\}$ , and  $y$  by  $y = \lim_{\varepsilon_{2j'}} x$ . This process must terminate after at most  $t$  steps.  $\square$

*Algorithm 16.14. [Covering by Contractible Sets]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $R$ .
- **Input:**
  - a finite set of  $s$  polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$  in strong  $k$ -general position on  $R^k$ , with  $\deg(P_i) \leq d$  for  $1 \leq i \leq s$ ,
  - a  $\mathcal{P}$ -closed semi-algebraic set  $S$ , contained in the sphere of center 0 and radius  $r$ , defined by a  $\mathcal{P}$ -closed formula  $\phi$ .
- **Output:** a set of formulas  $\{\phi_1, \dots, \phi_M\}$  defined by  $\wedge$ -polynomials in  $D_1[X_1, \dots, X_k]$  such that
  - each  $\text{Reali}(\phi_i, R_1^k)$  is semi-algebraically contractible, and
  - $\bigcup_{1 \leq i \leq M} \text{Reali}(\phi_i, R_1^k) = \text{Ext}(S, R_1)$ .
- **Complexity:** The complexity of the algorithm is bounded by  $s^{(k+1)^2} d^{O(k^5)}$ .
- **Procedure:**
  - Step 1: Let  $Q = X_1^2 + \dots + X_{k+1}^2 - r^2$ . Call Algorithm 16.5 (Parametrized Bounded Connecting) with input  $Q, \mathcal{P}$ . Let  $\mathcal{A}$  be the family of polynomials output.
  - Step 2: Compute the set of realizable sign conditions  $\text{SIGN}(S)$  using Algorithm 13.1 (Computing Realizable Sign Conditions) .
  - Step 3: Using Algorithm 14.21 (Quantifier Elimination), eliminate one variable to compute the image of the semi-algebraic map  $\gamma_{\sigma^-}$ . Finally, output the set of formulas  $\{\phi_\sigma \mid \sigma \in \text{SIGN}(\mathcal{A}, S)\}$  describing the semi-algebraic set  $C(\sigma)$ .

**Proof of correctness:** The correctness of the algorithm is a consequence of Proposition 16.16, Proposition 16.18 and the correctness of Algorithm 16.5 (Parametrized Bounded Connecting), as well as the correctness of Algorithm 13.1 (Computing Realizable Sign Conditions) and Algorithm 14.5 (Quantifier Elimination).  $\square$

**Complexity analysis:** The complexity of Step 1 of the algorithm is bounded by  $s^{k+1}d^{O(k^4)}$ , where  $s$  is a bound on the number of elements of  $\mathcal{P}$  and  $d$  is a bound on the degrees of the elements of  $\mathcal{P}$ , using the complexity analysis of Algorithm 16.5 (Parametrized Bounded Connecting). The number of polynomials in  $\mathcal{A}$  is  $s^{k+1}d^{O(k^4)}$  and their degrees are bounded by  $d^{O(k^3)}$ . Thus the complexity of computing  $\text{SIGN}(S)$  is bounded by  $s^{(k+1)^2}d^{O(k^5)}$  using Algorithm 13.1 (Computing Realizable Sign Conditions). In Step 3 of the algorithm there is a call to Algorithm 14.5 (Quantifier Elimination). There are two blocks of variables of size  $k$  and 2 respectively. The number and degrees of the input polynomials are bounded by  $s^{k+1}d^{O(k^4)}$  and  $d^{O(k^3)}$  respectively. Moreover, observe that even though we introduced  $2s$  infinitesimals, each arithmetic operation is performed in the ring  $D$  adjoined with at most  $O(k)$  infinitesimals since the polynomials  $\{P, P \pm \varepsilon_{2j}, P \pm \varepsilon_{2j-1}, P \in \mathcal{P}, 1 \leq j \leq s\}$  are in strong general position. Thus, the complexity of this step is bounded by  $s^{(k+1)^2}d^{O(k^5)}$  using the complexity analysis of Algorithm 14.5 (Quantifier Elimination) and the fact that each arithmetic operation costs at most  $d^{O(k^5)}$  in terms of arithmetic operations in the ring  $D$ .  $\square$

We now want to compute the first Betti number of a  $\mathcal{P}$ -closed semi-algebraic set  $S$  when  $\mathcal{P}$  is not necessary in general position. We first replace  $S$  by a  $\mathcal{P}^*$ -closed and bounded semi-algebraic set, where the elements of  $\mathcal{P}^*$  are slight modifications of the elements of  $\mathcal{P}$ , and the family  $\mathcal{P}^*$  is in general position and  $b_i(S^*) = b_i(S), 0 \leq i \leq k$ .

Define

$$H_i = 1 + \sum_{1 \leq j \leq k} i^j X_j^{d'}$$

where  $d'$  is the smallest number strictly bigger than the degree of all the polynomials in  $\mathcal{P}$ . Using arguments similar to the proof of Proposition 13.6 it is easy to see that the family  $\mathcal{P}^*$  of polynomials  $P_i - \delta H_i, P_i + \delta H_i$ , with  $P_i \in \mathcal{P}$ . is in general position in  $\mathbb{R}\langle \delta \rangle^k$ .

**Lemma 16.20.** *Denote by  $S^*$  the set obtained by replacing any  $P_i \geq 0$  in the definition of  $S$  by  $P_i \geq -\delta H_i$  and every  $P_i \leq 0$  in the definition of  $S$  by  $P_i \leq \delta H_i$ . If  $S$  is bounded, the set  $\text{Ext}(S, \mathbb{R}\langle \delta \rangle^k$  is semi-algebraically homotopy equivalent to  $S^*$ .*

**Proof:** The claim follows by Lemma 16.17. Note that  $S$  is closed and bounded,  $\lim_{\delta} S^* = S$ , and  $S^*(t) \subset S^*(t')$  for  $t < t'$ .  $\square$

*Algorithm 16.15. [First Betti Number in the  $\mathcal{P}$ -closed case]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$
- **Input:**
  - a finite set of polynomials  $\mathcal{P} \subset D[X_1, \dots, X_k]$ ,

- a formula defining a  $\mathcal{P}$ -closed semi-algebraic set,  $S$ .
- **Output:** the first Betti number  $b_1(S)$ .
- **Complexity:**  $(s d)^{k^{O(1)}}$ , where  $s = \#\mathcal{P}$  and  $d = \max_{P \in \mathcal{P}} \deg(P)$ .
- **Procedure:**
  - Step 1: Let  $\varepsilon$  be an infinitesimal. Replace  $S$  by the semi-algebraic set  $T$  defined as the intersection of the cylinder  $S \times \mathbb{R}\langle\varepsilon\rangle$  with the upper hemisphere defined by  $\varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) = 1, X_{k+1} \geq 0$ .
  - Step 2: Replace  $T$  by  $T^*$  using the notation of Lemma 16.20.
  - Step 3: Use Algorithm 16.14 (Covering by Contractible Sets) with input  $\varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 4$  and  $\mathcal{P}^*$ , to compute a covering of  $T^*$  by closed, bounded and contractible sets,  $T_i$ , described by formulas  $\phi_i$ .
  - Step 4: Use Algorithm 16.13 (General Roadmap) to compute exactly one sample point of each connected component of the pairwise and triplewise intersections of the  $T_i$ 's. For every pair  $i, j$  and every  $k$  compute the incidence relation between the connected components of  $T_{ijk}^*$  and  $T_{ij}^*$  as follows: compute a roadmap of  $T_{ij}^*$ , containing the sample points of the connected components of  $T_{ijk}^*$  using Algorithm 16.13 (General Roadmap).
  - Step 5: Using linear algebra compute

$$b_1(T^*) = \dim(\text{Ker}(\delta_2)) - \dim(\text{Im}(\delta_1)),$$

with

$$\prod_i H^0(T_i^*) \xrightarrow{\delta_1} \prod_{i < j} H^0(T_{ij}^*) \xrightarrow{\delta_2} \prod_{i < j < \ell} H^0(T_{ij\ell}^*)$$

**Proof of correctness:** First note that  $T$  is closed and bounded and has the same Betti numbers as  $S$ , using the local conical structure at infinity. It follows from Lemma 16.20 that  $T$  and  $T^*$  have the same Betti numbers. The correctness of the algorithm is a consequence of the correctness of Algorithm 16.14 (Covering by Contractible Sets), Algorithm 16.13 (General Roadmap), and Theorem 6.9. □

**Complexity analysis:** The complexity of Step 3 of the algorithm is bounded by  $s^{(k+1)^2} d^{O(k^6)}$  using the complexity analysis of Algorithm 16.14 (Covering by Contractible Sets) and noticing that each arithmetic operation takes place a ring consisting of  $D$  adjoined with at most  $k$  infinitesimals. Finally, the complexity of Step 4 is also bounded by  $(s d)^{k^{O(1)}}$ , using the complexity analysis of Algorithm 16.13 (General Roadmap). □

Now we describe the algorithm for computing the first Betti number of a general semi-algebraic set. We first replace the given set by a closed and bounded one, using Theorem 7.45. We then apply Algorithm 16.15.

*Algorithm 16.16. [First Betti Number of a  $\mathcal{P}$ -Semi-algebraic Set]*

- **Structure:** an ordered domain  $D$  contained in a real closed field  $\mathbb{R}$ .

- **Input:**
  - a finite set of polynomials  $\mathcal{P} \subset \mathbb{D}[X_1, \dots, X_k]$ ,
  - a formula defining a  $\mathcal{P}$ -semi-algebraic set,  $S$ .
- **Output:** the first Betti number  $b_1(T)$ .
- **Complexity:**  $(s d)^{k^{O(1)}}$ , where  $s = \#\mathcal{P}$  and  $d = \max_{P \in \mathcal{P}} \deg(P)$ .
- **Procedure:**
  - Step 1: Let  $\varepsilon$  be an infinitesimal. Define  $\tilde{S}$  as the intersection of  $\text{Ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  with the ball of center 0 and radius  $1/\varepsilon$ . Define  $\mathcal{Q}$  as
 
$$\mathcal{P} \cup \{\varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) - 4, X_{k+1}\}.$$
  - Replace  $\tilde{S}$  by the  $\mathcal{Q}$ -semi-algebraic set  $S$  defined as the intersection of the cylinder  $\tilde{S} \times \mathbb{R}\langle\varepsilon\rangle$  with the upper hemisphere defined by  $\varepsilon^2(X_1^2 + \dots + X_k^2 + X_{k+1}^2) = 4, X_{k+1} \geq 0$ .
  - Step 2: Using Definition 7.44, replace  $T$  by a  $\mathcal{Q}'$ -closed set,  $T'$ , where
 
$$\mathcal{Q}' = \{P \pm \varepsilon_i \mid P \in \mathcal{Q}, i = 1, \dots, 2s\}.$$
  - Step 3: Use Algorithm 16.15 to compute the first Betti number of  $T'$ .

**Proof of correctness:** The correctness of the algorithm is a consequence of Theorem 7.45 and the correctness of Algorithm 16.15. □

**Complexity analysis:** In Step 2 of the algorithm the cardinality of  $\mathcal{Q}'$  is  $2(s+1)^2$  and the degrees of the polynomials in  $\mathcal{Q}'$  are still bounded by  $d$ . The complexity of Step 3 of the algorithm is then bounded by  $(s d)^{k^{O(1)}}$  using the complexity analysis of Algorithm 16.15. □

## 16.7 Bibliographical Notes

A motivation for deciding connectivity of semi-algebraic sets comes from robot motion planning [146]. This is equivalent to deciding whether the two corresponding points in the free space are in the same connected component of the free space. The solution by Schwartz and Sharir [146] using Collin’s method of cylindrical algebraic decomposition. The complexity of their solution is thus polynomial in  $d$  and  $s$  and doubly exponential in  $k$ .

Canny introduced the notion of a roadmap for a semi-algebraic set and gave an algorithm [36] which after subsequent modifications [38] constructed a roadmap for a semi-algebraic set defined by polynomials whose sign invariant sets give a stratification of  $\mathbb{R}^k$  and whose complexity is  $s^k (\log s) d^{O(k^4)}$ . For an arbitrary semi-algebraic set he perturbs the defining polynomials and is then able to decide if two points are in the same semi-algebraically connected component with the same complexity. However, this algorithm does not give a path joining the points. A Monte Carlo version of this algorithm has complexity  $s^k (\log s) d^{O(k^2)}$ .

Grigor'ev and Vorobjov [78] gave an algorithm with complexity  $(sd)^{k^{O(1)}}$ , for counting the number of connected components of a semi-algebraic set. Heintz, Roy, and Solerno [86] and Gournay and Risler [75] gave algorithms which compute a roadmap for any semi-algebraic set whose complexity was also  $(sd)^{k^{O(1)}}$ .

Unlike the complexity of Canny's algorithm, the complexities of these algorithms are not separated into a combinatorial part (the part depending on  $s$ ) and an algebraic part (the part depending on  $d$ ). Since the given semi-algebraic set might have  $(sd)^k$  different connected components, the combinatorial complexity of Canny's algorithm is nearly optimal. Canny's algorithm makes use of Thom's isotopy lemma for stratified sets and consequently requires the use of generic projections, as well as perturbations to put the input polynomials into general position in a very strong sense. In order to do this in a deterministic fashion,  $O(s + k^2)$  different transcendentals are introduced, requiring the algebraic operations to be performed over an extended ring. This raises the algebraic complexity of the deterministic algorithm to  $d^{O(k^4)}$ .

In [16] a deterministic algorithm constructing a roadmap for any semi-algebraic set contained in an algebraic set  $\text{Zer}(Q, \mathbb{R}^k)$  of dimension  $k'$  with complexity  $s^{k'+1}d^{O(k^2)}$  is given. In robot motion planning, the configuration space of a robot is often embedded as a lower dimensional algebraic set in a higher dimensional real Euclidean space (see [103]), so it is of interest to design algorithms which take advantage of this fact and whose complexity reflects the dimension of this algebraic set rather than the dimension of the ambient space. The combinatorial complexity of this algorithm is nearly optimal. The algorithm uses only a fixed number of infinitesimal quantities which reduces the algebraic complexity to  $d^{O(k^2)}$ . The algorithm also computes a semi-algebraic path between the input points if they happen to lie in the same connected component and hence solves the full version of the problem.

A single exponential bound  $(sd)^{k^{O(1)}}$  for computing the connected components of a semi-algebraic set is due to Canny, Grigor'ev, Vorobjov and Heintz, Roy and Solerno [39, 87]. The results presented here are significantly more precise.

---

## References

1. J. ABDELJAOUED, H. LOMBARDI, *Méthodes matricielles: Introduction à la Complexité Algébrique*, Mathématiques & Applications, Springer (2004).
2. J.W. ALEXANDER, *A proof of the invariance of certain constants of analysis situs*, Trans. Amer. Math. Soc., 16 148–154 (1915).
3. N. ALON, *Tools from higher algebra*, Handbook of combinatorics, 1749–1783, Elsevier, Amsterdam (1995).
4. M. E. ALONSO, E. BECKER, M.-F. ROY, T. WÖRMANN, *Zeroes, Multiplicities and Idempotents for Zerodimensional Systems*, Algorithms in algebraic geometry and applications, Progress in Mathematics, vol 143. Birkhauser 1–16 (1996).
5. C. ANDRADAS, L. BRÖCKER, J. RUIZ, *Constructible sets in Real Geometry*, Ergeb. Math. Grenzgeb. (3) 33 Berlin etc.: Springer-Verlag (1996).
6. E. ARTIN, *Über die Zerlegung definiter Funktionen in Quadrate*, Hamb. Abh. 5, 100–115 (1927). The collected papers of Emil Artin, 273–288. Reading: Addison-Wesley (1965).
7. E. ARTIN, O. SCHREIER, *Algebraische Konstruktion reeller Körper*, Hamb. Abh. 5 8(-99) (1925). The collected papers of Emil Artin, 258–271. Addison-Wesley (1965).
8. P. AUBRY, F. ROUILLIER, M. SAFEY EL DIN, *Real solving for positive dimensional systems*, Journal of Symbolic Computation, 34(6), 543–560 (2002).
9. E. H. BAREISS, *Sylvester’s Identity and Multistep Integer-Preserving Gaussian Elimination*, Math. Comp. 22 565–578 (1968).
10. S. BASU, *On different bounds on different Betti numbers*, Discrete and Computational Geometry, 30:1, 65–85, (2003).
11. S. BASU, *On bounding the Betti Numbers and Computing the Euler Characteristics of Semi-algebraic Sets*, Discrete and Computational Geometry, 22 1–18 (1999).
12. S. BASU, *New Results on Quantifier Elimination over Real Closed Fields and Applications to Constraint Databases*, Journal of the ACM, 46 (4), 537–555 (1999).
13. S. BASU, R. POLLACK, M.-F. ROY, *On the Combinatorial and Algebraic Complexity of Quantifier Elimination*, Journal of the ACM, 43 1002–1045, (1996).
14. S. BASU, R. POLLACK, M.-F. ROY, *On the number of cells defined by a family of polynomials on a variety*, Mathematika, 43 120–126 (1996).

15. S. BASU, R. POLLACK, M.-F. ROY, *On Computing a Set of Points meeting every Semi-algebraically Connected Component of a Family of Polynomials on a Variety*, Journal of Complexity, March 1997, 13 (1), 28–37.
16. S. BASU, R. POLLACK, M.-F. ROY, *Computing Roadmaps of Semi-algebraic Sets on a Variety*, Journal of the AMS, 3 (1) 55–82 (1999).
17. S. BASU, R. POLLACK, M.-F. ROY *On the Betti Numbers of Sign Conditions*, Proc. Amer. Math. Soc. 133, 965–974 (2005).
18. S. BASU, R. POLLACK, M.-F. ROY, *Computing Euler-Poincaré characteristic of sign conditions*, Journal of Computational Complexity, 14, 53–71 (2005).
19. S. BASU, R. POLLACK, M.-F. ROY, *Computing the Dimension of a Semi-Algebraic Set*, Zap. Nauchn. Semin. POMI 316 42–54 (2004).
20. S. BASU, R. POLLACK, M.-F. ROY, *Algorithms in real algebraic geometry*, Springer-Verlag (2003).
21. S. BASU, R. POLLACK, M.-F. ROY, *Computing the first Betti number and describing the connected components of semi-algebraic sets*, preprint (2005), [arXiv:math.AG/0603248].
22. R. BENEDETTI, J.-J. RISLER, *Real algebraic and semi-algebraic sets*, Actualités Mathématiques. Hermann, Paris (1990).
23. M. BEN-OR, D. KOZEN, J. REIF, *The complexity of elementary algebra and geometry*, J. of Computer and Systems Sciences, 18:251–264, (1986).
24. E. BÉZOUT, *Recherche sur les degrés de l'équation résultant de l'évanouissement des inconnues*, Histoire de l'académie royale des sciences, 288–338, (1764).
25. G.D. BIRKHOFF, *Collected Mathematical Papers*, Vol II, Amer. Math. Soc., New York (1950).
26. J. BOCHNAK, M. COSTE, M.-F. ROY, *Géométrie algébrique réelle*, Springer-Verlag (1987). *Real algebraic geometry*, Springer-Verlag (1998).
27. A. BOREL, J. C. MOORE, *Homology theory for locally compact spaces*, Mich. Math. J., 7: 137–159, (1960).
28. H. BRAKHAGE, *Topologische Eigenschaften algebraischer Gebilde über einen beliebigen reell-abgeschlossenen Konstantenkörper*, Dissertation, Univ. Heidelberg (1954).
29. E. BRIAND, F. CARRERAS, L. GONZALEZ-VEGA, N. GONZALEZ-CAMPOS, I. NECULA, H. PERDRY, N. DEL RIO, C. TANASESCU, F. ROUILLIER, M.-F. ROY, A. SEIDL, *Creating an Electronic Book for Algorithms in Real Algebraic Geometry: a first experiment*, ISSAC Poster (2004).  
<http://www0.risc.uni-linz.ac.at/issac2004/poster-abstracts/abstract24.pdf>
30. C. BROWN, *QEPCAD B: a program for computing with semi-algebraic sets using CADs* ACM SIGSAM 37 (4): 97–108 (2003).  
<http://www.cs.usna.edu/~qepcad/B/QEPCAD.html>
31. W. D. BROWNAWELL *Local diophantine Nullstellen identities*, J. AMS 1, 311–322 (1988).
32. B. BUCHBERGER, *An Algorithm for Finding the Basis Elements in the Residue Class Ring Modulo a Zero Dimensional Polynomial Ideal*, PhD Thesis, Mathematical Institute, University of Innsbruck, Austria, (1965).
33. B. BUCHBERGER, *Gröbner bases: an algorithmic method in polynomial ideal theory*, Recent trends in multidimensional systems theory, Reider ed. Bose, (1985).

34. F. BUDAN DE BOISLAURENT, *Nouvelle méthode pour la résolution des équations numériques d'un degré quelconque*, (1807), 2nd edition, Paris (1822).
35. L. CANIGLIA, A. GALLIGO, J. HEINTZ, *Borne simplement exponentielle pour les degrés dans le théorème des zéros sur un corps de caractéristique quelconque*, C. R. Acad. Sci. Paris 307, 255–258 (1988).
36. J. CANNY, *The Complexity of Robot Motion Planning*, MIT Press (1987).
37. J. CANNY, *Some Algebraic and Geometric Computations in PSPACE*, Proc. Twentieth ACM Symp. on Theory of Computing, 460–467, (1988).
38. J. CANNY, *Computing road maps in general semi-algebraic sets*, The Computer Journal, 36: 504–514, (1993).
39. J. CANNY, D. GRIGOR'EV, N. VOROBJOV, *Finding connected components of a semi-algebraic set in subexponential time*, Appl. Algebra Eng. Commun. Comput., 2 (4), 217–238 (1992).
40. F. CARUSO, *SARAG: Some Algorithms in Real Algebraic Geometry*, (2005).
41. A. CAUCHY, *Calcul des indices des fonctions*, Journal de l'Ecole Polytechnique, vol. 15, Cahier 25, 176–229 (1832).
42. A. CHISTOV, H. FOURNIER, L. GURVITS, P. KOIRAN, *Vandermonde Matrices, NP-Completeness and Transversal Subspaces*, Foundations of Computational Mathematics, 3 (4) 421–427 (2003).
43. P. J. COHEN, *Decision procedures for real and  $p$ -adic fields*, Comm. Pure. Appl. Math. 22, 131–151 (1969).
44. G. E. COLLINS, *Subresultants and Reduced Polynomial Remainder Sequences*, Journal of the ACM 14 128–142 (1967).
45. G. COLLINS, *Quantifier elimination for real closed fields by cylindric algebraic decomposition*, In Second GI Conference on Automata Theory and Formal Languages. Lecture Notes in Computer Science, vol. 33, 134–183, Springer-Verlag, Berlin (1975).
46. G. E. COLLINS, H. HONG, *Partial Cylindrical Algebraic Decomposition for Quantifier Elimination*, Journal of Symbolic Computation, 12, 299–328 (1991).
47. M. COSTE, *An introduction to semi-algebraic geometry*, Dip. Mat. Univ. Pisa, Dottorato di Ricerca in Matematica, Istituti Editoriali e Poligrafici Internazionali, Pisa (2000).
48. M. COSTE, *An introduction to  $o$ -minimal geometry*, Dip. Mat. Univ. Pisa, Dottorato di Ricerca in Matematica, Istituti Editoriali e Poligrafici Internazionali, Pisa (2000).
49. M. COSTE, T. LAJOUS-LOEZA, H. LOMBARDI, M.-F. ROY, *Generalized Budan-Fourier theorem and virtual roots*, Journal of Complexity, 21, 478–486, (2005).
50. M. COSTE, M.-F. ROY, *Thom's lemma, the coding of real algebraic numbers and the topology of semi-algebraic sets*. Journal of Symbolic Computation 5, No.1/2, 121–129 (1988).
51. D. COX, J. LITTLE, D. O'SHEA, *Ideals, varieties and algorithms: an introduction to computational algebraic geometry and commutative algebra*, Undergraduate Texts in Mathematics, Springer-Verlag, New York (1997).
52. J. H. DAVENPORT, J. HEINTZ, *Real quantifier elimination is doubly exponential*, Journal of Symbolic Computation 5, No.1/2, 29–35 (1988).
53. R. DESCARTES, *Géométrie* (1636). A source book in Mathematics, 90–31. Harvard University press (1969).



54. C. L. DOGDSON, *Condensation of determinants, being a new and brief method for computing their numerical values*, Proc. Royal. Soc. Lond. 15 150–155 (1866).
55. A. EIGENWILLIG, V. SHARMA, C. YAP, *Almost tight complexity bounds for Descartes method*, preprint (2006).
56. L. EULER, *Démonstration sur le nombre de points où deux lignes d'ordre quelconque peuvent se couper*, Mémoires de l'Académie des Sciences de Berlin, 4, 234–248, (1750).
57. G. FARIN, *Curves and surfaces for Computer Aided Design*, Academic Press (1990).
58. J.-C. FAUGÈRE, *A new efficient algorithm for computing Gröbner basis (F 4)* Journal of Pure and Applied Algebra, 139 61–88 (1999).
59. J.-C. FAUGÈRE, *FGb (Gröbner basis computations)*  
<http://fgbrs.lip6.fr/Software/>
60. J. FOURIER, *Analyse des équations déterminées*, F. Didot, Paris (1831).
61. F. G. FROBENIUS, *Über das Traegheitsgesetz des quadratischen Formen*, S-B Pruss. Akad. Wiss. 241-256 (1884), 403–431 (1884).
62. A. GABRIELOV, N. VOROBOV, *Betti Numbers for Quantifier-free Formulae*. Discrete and Computational Geometry, 33 395–401 (2005).
63. F. R. GANTMACHER, *Theory of matrices, Vol I*. AMS-Chelsea (2000).
64. J. VON ZUR GATHEN, J. GERHARD *Modern computer algebra*, Cambridge University Press (1999).
65. C. F. GAUSS *Demonstratio Nova Altera Theorematis Omnem Funct. Alg.*, Commentationes societatis regieae scientiarum Gottingensis recentiores, 3, 107–134 (1816). Werke III 31-56 (1876).
66. L. GONZALEZ VEGA, *La sucesión de Sturm-Habicht y sus aplicaciones al Algebra Computacional*, Doctoral Thesis, Universidad de Cantabria (1989).
67. L. GONZALEZ VEGA, M. EL KAHOUI, *An improved upper complexity bound for the topology computation of a real algebraic curve*, J. Complexity 12 527–544 (1996).
68. L. GONZALEZ VEGA, H. LOMBARDI, L. MAHÉ, *Virtual roots of real polynomials* Journal of Pure and Applied Algebra, 124, 147–166 (1998).
69. L. GONZALEZ VEGA, H. LOMBARDI, T. RECIO, M.-F. ROY, *Spécialisation de la suite de Sturm et sous-résultants I*, Informatique théorique et applications 24 561–588 (1990).
70. L. GONZALEZ VEGA, H. LOMBARDI, T. RECIO, M.-F. ROY, *Spécialisation de la suite de Sturm et sous-résultants II*, Informatique théorique et applications 28 1-24 (1994).
71. L. GONZALEZ-VEGA, I. NECULA, *Efficient topology determination of implicitly defined algebraic plane curves*, Comput. Aided Geom. Design 19, no. 9, 719–743 (2004).
72. L. GONZALEZ VEGA, F. ROUILLIER, M.-F. ROY, *Symbolic Recipes for Polynomial System Solving*, In: Some tapas of computer algebra, A. Cohen et al. ed. Algorithms and Computation in Mathematics, vol. 4, 34–64, Springer.
73. L. GONZALEZ VEGA, F. ROUILLIER, M.-F. ROY, G. TRUJILLO *Symbolic Recipes for Real Solutions*, In: Some tapas of computer algebra, A. Cohen et al. ed. Algorithms and Computation in Mathematics, vol. 4, 121–167, Springer.
74. P. GORDAN, *Verlesungen über Invariantentheorie- Ester Band: Determinanten*, B. G. Teubner, Leipzig (1885).

75. L. GOURNAY, J. J. RISLER, *Construction of roadmaps of semi-algebraic sets*, Appl. Algebra Eng. Commun. Comput. 4, No.4, 239–252 (1993).
76. D. GRIGOR'EV, *The Complexity of deciding Tarski algebra*, Journal of Symbolic Computation 5 65–108 (1988).
77. D. GRIGOR'EV, N. VOROBYOV, *Solving Systems of Polynomial Inequalities in Subexponential Time*, Journal of Symbolic Computation, 5 37–64 (1988).
78. D. GRIGOR'EV, N. VOROBYOV, *Counting connected components of a semi-algebraic set in subexponential time*, Comput. Complexity 2, No.2, 133–186 (1992).
79. W. HABICHT, *Eine Verallgemeinerung des Sturmschen Wurzelzählverfahrens*, Comm. Math. Helvetici 21, 99–116 (1948).
80. J. HADAMARD, *Résolution d'une question relative au déterminant*, Bull. Sci. Math. 17, 240–246 (1893).
81. HAM, LE, *Un theoreme de Zariski du type de Lefschetz*, Ann. Sc. Ec. Norm. Sup., (3) vol 6. 317–355 (1973).
82. HAM, LE, *Lefschetz theorems on quasi-projective varieties*, Bull. Soc. Math. France, 113, 123–142 (1985).
83. R. M. HARDT, *Semi-algebraic Local Triviality in Semi-algebraic Mappings*, Am. J. Math. 102, 291–302 (1980).
84. J. HEINTZ, M.-F. ROY, P. SOLERNÒ, *On the Complexity of Semi-Algebraic Sets*, Proc. IFIP 89, San Francisco. North- Holland 293–298 (1989).
85. J. HEINTZ, M.-F. ROY, P. SOLERNÒ, *Sur la complexité du principe de Tarski-Seidenberg*, Bull. Soc. Math. France 118 101–126 (1990).
86. J. HEINTZ, M.-F. ROY, P. SOLERNÒ, *Single exponential path finding in semi-algebraic sets II: The general case*, Bajaj, Chandrajit L. (ed.), Algebraic geometry and its applications. Collections of papers from Shreeram S. Abhyankar's 60th birthday conference held at Purdue University, West Lafayette, IN, USA, June 1-4, 1990. New York: Springer-Verlag, 449–465 (1994).
87. J. HEINTZ, M.-F. ROY, P. SOLERNÒ, *Description of the Connected Components of a Semialgebraic Set in Single Exponential Time*, Discrete and Computational Geometry, 11, 121–140 (1994).
88. G. HERMANN, *Die Frage der endlich vielen Schritte in der Theorie der Polynomideale*, Math. Annalen 95, 736–788 (1926).
89. C. HERMITE, *Remarques sur le théorème de Sturm*, C. R. Acad. Sci. Paris 36, 52–54 (1853).
90. C. HERMITE, *Sur l'indice des fractions rationnelles*, Bull. Soc. Math. France, tome 7, 128–131 (1879).
91. D. HILBERT, *Über die Theorie der algebraischen Formen*, Math. Annalen 36, 473–534 (1890).
92. L. HÖRMANDER, *The analysis of linear partial differential operators*, vol. 2. Berlin etc.: Springer-Verlag (1983).
93. A. HURWITZ, *Über die Bedingungen, unter welchen eine Gleichung nur Wurzeln mit negativen reellen Theilen besitzt*, Math. Annalen, vol. 46, 273–284 (1895).
94. N. V. ILYUSHECKIN, *On some identities for the elements of a symmetric matrix*, Zapiski Nauchnyh Seminarov POMI Vol. 303, Investigations on Linear Operators and Function Theory. Part 31, editor S.V.Kislyakov (1997), <http://www.pdmi.ras.ru/zns1/2003/v303.html>.

95. A. KHOVANSKY, *Fewnomials*, Transl. Math. Monogr. 88, Providence, RI: American Mathematical Society (1991).
96. M. KNEBUSCH, C. SCHEIDERER, *Einführung in die reelle Algebra*, Vieweg-Studium 63, Aufbaukurs Mathematik, Vieweg (1989).
97. J. KOLLAR, *Effective Nullstellensatz for arbitrary ideals*, J. Eur. Math. Soc. 1 no. 3, 313–337 (1999).
98. W. KRANDICK, K. MEHLHORN, *New bounds for the Descartes method*, Journal of Symbolic Computation, vol. 41, 49–66 (2006).
99. M. G. KREIN AND M. A. NAIMARK, *The method of symmetric and hermitian forms in the theory of separation of roots of algebraic equations*, Kharkov 1936 (in Russian), English translation in: Linear and Multi-linear Algebra 10 265–308 (1981).
100. L. KRONECKER, *Werke*, Vol 1. Leipzig, Teubner (1895).
101. J.-M. LANE, R. F. RIESENFELD, *Bounds on a polynomial*, 21 112–117 (1981).
102. S. LANG, *Algebra*, Reading: Addison-Wesley (1971).
103. J.-C. LATOMBE, *Robot Motion Planning*, The Kluwer International Series in Engineering and Computer Science. 124. Dordrecht: Kluwer Academic Publishers Group (1991).
104. D. LAZARD, *Gröbner-Bases, Gaussian elimination and resolution of systems of algebraic equations*. EUROCAL 146–156 (1983).
105. S. LEFSCHETZ, *Algebraic geometry*. Princeton University Press, Princeton, N. J., (1953).
106. U. LE VERRIER *Sur les variations séculaires des éléments elliptiques des sept planètes principales: Mercure, Vénus, La Terre, Mars, Jupiter, Saturne et Uranus*, J. Math. Pures Appli., (4), 220–254 (1840).
107. T. LICKTEIG, M.-F. ROY, *Sylvester-Habicht sequences and fast Cauchy index computation*, Journal of Symbolic Computation, 31 315–341 (2001).
108. A. LIÉNARD, M. H. CHIPART *Sur le signe de la partie réelle des racines d’une équation algébrique*, J. Math. Pures Appl. (6) 10 291–346 (1914).
109. Z. LIGATSIKAS, M.-F. ROY, *Séries de Puiseux sur un corps réel clos*. C. R. Acad. Sci. Paris 311 625–628 (1990).
110. S. LOJASIEWICZ *Ensembles semi-analytiques*. Inst. Hautes Etudes Sci., (preprint) (1964).
111. S. LOJASIEWICZ *Triangulation of semi-analytic sets*. Ann. Scuola Norm. Sup. Pisa, Sci. Fis. Mat. (3) 18, 449–474 (1964).
112. H. LOMBARDI, M.-F. ROY, M. SAFEY, *New structure theorems for subresultants*, Special Issue Symbolic Computation in Algebra, Analysis, and Geometry Journal of Symbolic Computation, 29 663–690 (2000).
113. Y. MATIYASEVICH *Hilbert’s tenth problem*. Translated from the 1993 Russian original by the author. With a foreword by Martin Davis. Foundations of Computing Series. MIT Press, Cambridge, MA (1993).
114. E. W. MAYR *Some complexity results for polynomial ideals*, Journal of Complexity, 13(3):303–325 (1997).
115. E. MENDELSON *Introduction to mathematical logic*, Princeton, N.J., Van Nostrand (1964).
116. M. MIGNOTTE, D. STEFANESCU *Polynomials, an algorithmic approach*, Springer Verlag, Singapore (1999).
117. J. MILNOR *Morse Theory*, Annals of Mathematical Studies, Princeton University Press (1963).

118. J. MILNOR, *On the Betti numbers of real varieties*, Proc. AMS 15, 275–280 (1964).
119. R. T. MOENCK *Fast computation of GCDs*, Proc. STOC '73, 142–151 (1973).
120. L. G. MONK, *Elementary-recursive decision procedures*, Thesis (1976).
121. M. MORSE, *Relations between the critical points of a real function of  $n$  independent variables*, Trans. Amer. Math. Soc., 27 345–396 (1925).
122. B. MOURRAIN, M. N. VRAHATIS, J.-C. YAKHOUBSON *On the Complexity of Isolating Real Roots and Computing with Certainty the Topological Degree*, Journal of Complexity, 182, 612–640 (2002).
123. I. NEWTON, *The mathematical papers of Isaac Newton*, Cambridge University Press (1968,1971,1976).
124. O. A. OLEINIK, *Estimates of the Betti numbers of real algebraic hypersurfaces*, Mat. Sb. (N.S.), 28 (70): 635–640 (Russian) (1951).
125. O. A. OLEINIK, I. B. PETROVSKII, *On the topology of real algebraic surfaces*, Izv. Akad. Nauk SSSR 13, 389–402 (1949).
126. A. OSTROWSKI, *Notes sur les produits de séries normales*, Bulletin de la Société Royale des Sciences de Liège 8 458–467 (1939).
127. C. PAPADIMITRIOU, *Computational Complexity*, Addison-Wesley (1994).
128. P. PEDERSEN, *Counting real zeroes of polynomials*, PhD Thesis, Courant Institute, New York University (1991).
129. P. PEDERSEN, M.-F. ROY, A. SZPIRGLAS, *Counting real zeroes in the multivariate case*, Computational algebraic geometry, Eyssette et Galligo ed. Progress in Mathematics 109, 203–224, Birkhauser (1993).
130. H. POINCARÉ, *Analysis Situs*, Oeuvres, vol VI, Gauthier-Villars, Paris, 193–288 (1953).
131. R. POLLACK, M.-F. ROY, *On the number of cells defined by a set of polynomials*, C. R. Acad. Sci. Paris 316 573–577 (1993).
132. F. PREPARATA, D. SARWATE *An improved parallel processor bound in fast matrix inversion*, Inf. Proc. Letters, (7)/3, 148–150 (1978).
133. J. RENEGAR, *On the computational complexity and geometry of the first order theory of the reals*, Journal of Symbolic Computation, 13: 255–352 (1992).
134. F. ROUILLIER, *Solving Polynomial Systems through the Rational Univariate Representation*, Applicable Algebra in Engineering Communication and Computing, 9 (5) 433–461 (1999).
135. F. ROUILLIER, *RS (Real roots of systems with a finite number of complex solutions)*, <http://fgbrs.lip6.fr/Software/>
136. F. ROUILLIER, *On solving zero-dimensional polynomial systems with rational coefficients*, preprint (2005).
137. F. ROUILLIER, M.-F. ROY, M. SAFEY EL DIN, *Finding at least one point in each connected component of a real algebraic set defined by a single equation*, Journal of Complexity 16, 716–750 (2000).
138. F. ROUILLIER, P. ZIMMERMANN, *Efficient Isolation of a Polynomial Real Roots*, Rapport de Recherche INRIA 4113 (2001).
139. M.-F. ROUTH, *Stability of a given state of motion*, London (1877), The advanced part of a treatise of the system of rigid bodies Dover, New York (1955).
140. M.-F. ROY, *Basic algorithms in real algebraic geometry: from Sturm theorem to the existential theory of reals*, Lectures on Real Geometry in memoriam of Mario Raimondo, de Gruyter Expositions in Mathematics, 1–67 (1996).

141. M.-F. ROY, A. SZPIRGLAS, *Complexity of the computations with real algebraic numbers*, Journal of Symbolic computation 10 39–51 (1990).
142. M.-F. ROY, N. VOROBYOV, *Computing the Complexification of a Semi-algebraic Set*, Math. Zeitschrift 239, 131–142 (2002).
143. M. SAFEY EL DIN, *RAG'Lib*. <http://fgbrs.lip6.fr/Software/>
144. M. SAFEY EL DIN, E. SCHOST, *Properness defects of projections and computation of one point in each connected component of a real algebraic set*, Discrete and Computational Geometry, 32 (3), 417–430 (2004).
145. A. SCHONHAGE, *Schnelle Berechnung von Kettenbruchentwicklungen*, Acta Informatica 1,139–144, (1971).
146. J. SCHWARTZ, M. SHARIR, *On the ‘piano movers’ problem II. General techniques for computing topological properties of real algebraic manifolds*, Adv. Appl. Math. 4, 298–351 (1983).
147. I.R. SHAFAREVITCH, *Basic algebraic geometry*, Springer (1974).
148. A. SEIDENBERG, *A new decision method for elementary algebra*, Annals of Mathematics, 60:365–374, (1954).
149. V. SHARMA, C. YAP, *Sharp Amortized Bounds for Descartes’ and de Casteljau’s Methods for Real Root Isolation*, preprint (2005).
150. E. H. SPANIER, *Algebraic Topology*, McGraw-Hill Book Company (1966).
151. A. STREZEBONSKI, *Solving Systems of Strict Polynomial Inequalities*, Journal of Symbolic Computation 29 (3) 471–480 (2000).
152. C. STURM, *Mémoire sur la résolution des équations numériques*. Inst. France Sc.Math. Phys.6 (1835).
153. J. J. SYLVESTER, *On a theory of syzygetic relations of two rational integral functions, comprising an application to the theory of Sturm’s function*. Trans. Roy. Soc. London (1853).
154. A. TARSKI, *Sur les ensembles définissables de nombres réels*, Fund. Math. 17, 210–239 (1931).
155. A. TARSKI, *The completeness of elementary algebra and geometry*, 1939, Preprint Institut Blaise Pascal, CNRS (1967).
156. A. TARSKI, *A Decision method for elementary algebra and geometry*, University of California Press (1951).
157. R. THOM, *Sur l’homologie des variétés algébriques réelles*, Differential and Combinatorial Topology, 255–265. Princeton University Press, Princeton (1965).
158. J.V. USPENSKY, *Theory of equations*, MacGraw Hill (1948).
159. B. L. VAN DER WAERDEN, *Modern Algebra, Volume II*, F. Ungar Publishing Co. (1950).
160. B. L. VAN DER WAERDEN, *Topologische Begründung des Kalküls der abzählenden Geometrie*, Math. Ann. 102, 337–362 (1929).
161. L. VIETORIS, *Über die Homologiegruppen der Vereinigung zweier Komplexe*, Monatsh. für Math. u. Phys., 37 159–162 (1930).
162. A.J.H. VINCENT, *Sur la résolution des équations numériques*, Journal de Mathématiques Pures et Appliquées, 341–372 (1836).
163. N. N. VOROBYOV. *Complexity of computing the dimension of a semi-algebraic set*. J. of Symbolic Comput., 27: 565–579 (1999).
164. R. J. WALKER, *Algebraic Curves*, Princeton University Press (1950).

165. H.E. WARREN, *Lower bounds for approximations by non-linear manifolds*, Trans. Amer. Math. Soc. 1333, 167–178, (1968).
166. A. ZAPLETAL, *Gröbner bases bibliography*.  
<http://www.ricam.oeaw.ac.at/Groebner-Bases-Bibliography/index.php>

---

# Index of Notation

## Chapter 1

$\text{Zer}(\mathcal{P}, C^k)$	zero set of $\mathcal{P}$ in $C^k$ . . . . .	11
$\text{Free}(\Phi)$	free variables of a formula $\Phi$ . . . . .	12
$\text{Real}(\Phi, C^k)$	C-realization of a formula $\Phi$ in $C^k$ . . . . .	13
$\text{deg}(P)$	degree of $P$ . . . . .	15
$\text{cof}_j(P)$	coefficient of $X^j$ in $P$ . . . . .	15
$\text{lcof}(P)$	leading coefficient . . . . .	15
$\text{Rem}(P, Q)$	remainder in euclidean division . . . . .	15
$\text{Quo}(P, Q)$	quotient in euclidean division . . . . .	15
$\text{gcd}(P, Q)$	greatest common divisor of $P$ and $Q$ . . . . .	15
$\text{lcm}(P, Q)$	least common multiple of $P$ and $Q$ . . . . .	15
$\text{SRemS}(P, Q)$	signed remainder sequence of $P$ and $Q$ . . . . .	16
$\text{SRemS}_i(P, Q)$	$i$ -th signed remainder of $P$ and $Q$ . . . . .	16
$\text{Ex}(P, Q)$	extended signed remainder sequence of $P$ and $Q$ . . . . .	17
$\text{SRemU}_i(P, Q)$	$i$ -th cofactor of $P$ . . . . .	17
$\text{SRemV}_i(P, Q)$	$i$ -th cofactor of $Q$ . . . . .	17
$\text{gcd}(\mathcal{P})$	greatest common divisor of a family of polynomials $\mathcal{P}$ . . . . .	19
$\text{PRem}(P, Q)$	signed pseudo-remainder of $P$ and $Q$ . . . . .	21
$\text{Tru}_i(Q)$	truncation of $Q$ at degree $i$ . . . . .	21
$\text{Tru}(Q)$	set of truncations of $Q$ . . . . .	22
$\text{TRems}(P, Q)$	tree of possible signed pseudo-remainder sequence . . . . .	22
$\text{deg}_X(Q) = i$	basic formula fixing the degree of $Q$ to $i$ . . . . .	23
$\text{posgcd}(\mathcal{P})$	set of possible greatest common divisors of $\mathcal{P}$ . . . . .	24
$\text{Ext}(S, C')$	extension of $S$ to $C'$ . . . . .	27

## Chapter 2

$P'$	derivative of a polynomial $P$ . . . . .	29
$\text{sign}(a)$	sign of an element $a$ in an ordered field . . . . .	31
$ a $	absolute value of an element $a$ in an ordered field . . . . .	32
$0_+$	only order on $F(\varepsilon)$ such that $\varepsilon$ is infinitesimal over $F$ . . . . .	32

$F^{(2)}$	squares of elements in $F$ . . . . .	33
$\sum F^{(2)}$	sums of squares of elements in $F$ . . . . .	33
$E_i$	$i$ -th elementary symmetric function . . . . .	35
$\mathcal{M}_k$	set of monomials in $k$ variables . . . . .	35
$<_{\text{lex}}$	lexicographical ordering . . . . .	35
$<_{\text{grlex}}$	graded lexicographical ordering . . . . .	36
$ z $	modulus of $z$ . . . . .	40
$\text{Reali}(\sigma)$	realization of the sign condition $\sigma$ . . . . .	42
$\text{Der}(P)$	list of derivatives of $P$ . . . . .	42
$\text{Var}(a)$	number of sign variations in a sequence $a$ . . . . .	44
$\text{Var}(P)$	number of sign variations in the coefficients of $P$ . . . . .	44
$\text{pos}(P)$	number of positive roots of $P$ . . . . .	44
$\text{Var}(\mathcal{P}; a)$	number of sign variations of $\mathcal{P}$ at $a$ . . . . .	44
$\text{num}(P; (a, b])$	number of roots of $P$ in $(a, b]$ . . . . .	45
$v(P, x)$	virtual multiplicity of $x$ with respect to $P$ . . . . .	50
$v(P; (a, b])$	number of virtual roots of $P$ in $(a, b]$ . . . . .	50
$\text{Ind}(Q/P; a, b)$	Cauchy index of $Q/P$ on $(a, b)$ . . . . .	53
$\text{TaQ}(Q, P; a, b)$	Tarski-query of $Q$ for $P$ on $(a, b)$ . . . . .	54
$\text{Zer}(\mathcal{P}, \mathbb{R}^k)$	set of zeros of $\mathcal{P}$ in $\mathbb{R}^k$ . . . . .	57
$\text{Free}(\Phi)$	free variables of a formula $\Phi$ . . . . .	58
$\text{Reali}(\Phi, \mathbb{R}^k)$	realization of a formula $\Phi$ in $\mathbb{R}^k$ . . . . .	59
$\text{Reali}(\sigma)$	realization of the sign condition $\sigma$ over $Z$ . . . . .	64
$\text{Mat}(A, \Sigma)$	matrix of signs of $Q^A$ on $\Sigma$ . . . . .	64
$M \otimes M'$	tensor product of $M$ and $M'$ . . . . .	65
$M_s$	$3^s \times 3^s$ matrix defined inductively . . . . .	66
$\text{Ext}(S, \mathbb{R}')$	extension of $S$ to $\mathbb{R}'$ . . . . .	72
$\text{K}[[\varepsilon]]$	formal power series in $\varepsilon$ with coefficients in $\text{K}$ . . . . .	74
$\text{K}((\varepsilon))$	Laurent series in $\varepsilon$ with coefficients in $\text{K}$ . . . . .	74
$\text{K}\langle\langle\varepsilon\rangle\rangle$	Puiseux series in $\varepsilon$ with coefficients in $\text{K}$ . . . . .	74
$o(a)$	order of $a$ . . . . .	75
$\text{In}(a)$	initial coefficient of $a$ . . . . .	75
$Q(P, E, X)$	characteristic polynomial of an edge $E$ of the Newton polygon of $P$ . . . . .	76
$\text{K}(\varepsilon)$	algebraic Puiseux series in $\varepsilon$ with coefficients in $\text{K}$ . . . . .	81
$\text{K}(\varepsilon)_b$	algebraic Puiseux series with non-negative order . . . . .	81
$\lim_\varepsilon$	homomorphism from $\text{K}(\varepsilon)_b$ to $\text{K}$ mapping $\sum_{i \in \mathbb{N}} a_i \varepsilon^{i/q}$ to $a_0$ . . . . .	81
<b>Chapter 3</b>		
$\ x\ $	euclidean norm of $x$ . . . . .	83
$B_k(x, r)$	open ball in $\mathbb{R}^k$ with center $x$ and radius $r$ . . . . .	83
$\bar{B}_k(x, r)$	closed ball in $\mathbb{R}^k$ with center $x$ and radius $r$ . . . . .	83



$S^{k-1}(x, r)$	sphere in $\mathbb{R}^k$ with center $x$ and radius $r$ . . . . .	83
$\bar{S}$	closure of $S$ . . . . .	84
$S^\circ$	interior of $S$ . . . . .	84
$g \circ \phi$	composition of semi-algebraic function $g$ with a germ of a semi-algebraic continuous function $\phi$ . . . . .	91
$f'(x_0)$	derivative of a function $f$ at $x_0$ . . . . .	94
$df(x_0)$	derivative of a mapping $f$ at $x_0$ . . . . .	95
$\mathcal{S}^\infty(U, \mathbb{R})$	ring of Nash functions on $U$ . . . . .	95
$\ F\ $	norm of a linear mapping $F: \mathbb{R}^k \rightarrow \mathbb{R}^p$ . . . . .	95
$T_x(M)$	tangent space to $M$ at $x$ . . . . .	97

**Chapter 4**

Disc( $P$ )	discriminant of a monic polynomial $P$ . . . . .	101
sDisc $_k(P)$	subdiscriminant of a monic polynomial $P$ . . . . .	102
$N_n$	$n$ -th Newton sum of a polynomial . . . . .	102
Newt $_k(P)$	matrix with entries the Newton sums of $P$ . . . . .	103
det( $M$ )	determinant of a square matrix $M$ . . . . .	103
$V(x_1, \dots, x_r)$	Vandermonde matrix of $x_1, \dots, x_r$ . . . . .	104
Syl( $P, Q$ )	Sylvester matrix of $P$ and $Q$ . . . . .	106
Res( $P, Q$ )	resultant of $P$ and $Q$ . . . . .	106
SyHa $_j(P, Q)$	$j$ -th Sylvester-Habicht matrix of $P$ and $Q$ . . . . .	110
sRes $_j(P, Q)$	$j$ -th signed subresultant coefficient of $P$ and $Q$ . . . . .	110
$\varepsilon_i$	signature of the permutation reversing the order of $i$ consecutive rows in a matrix, i.e. $\varepsilon_i = (-1)^{i(i-1)/2}$ . . . . .	111
PmV( $s$ )	generalization of the difference between the number of sign permanences and the number of sign variations in the sequence $s$ . . . . .	113
sRes( $P, Q$ )	sequence of sRes $_j(P, Q)$ . . . . .	113
$A^t$	transpose of a matrix $A$ . . . . .	119
Rank( $A$ )	rank of a matrix $A$ . . . . .	119
Rank( $\Phi$ )	rank of a quadratic form $\Phi$ . . . . .	119
Sign( $\Phi$ )	signature of a quadratic form $\Phi$ . . . . .	121
$u \cdot u'$	inner product of $u$ and $u'$ . . . . .	122
$\ u\ $	norm of $u$ . . . . .	122
Tr( $M$ )	trace of the matrix $M$ . . . . .	125
Her( $P, Q$ )	Hermite quadratic form in the univariate case . . . . .	128
$L_f$	map of multiplication by $f$ in the univariate case . . . . .	128
mod $I$	congruence modulo the ideal $I$ . . . . .	132
$A/I$	quotient ring of an ideal $I$ of $A$ . . . . .	132
$\sqrt{I}$	radical of $I$ . . . . .	132
Ideal( $\mathcal{P}, \mathbb{K}$ )	ideal generated by $\mathcal{P}$ in $\mathbb{K}[X_1, \dots, X_k]$ . . . . .	132
$\mathcal{M}_k$	set of monomials in $k$ variables . . . . .	132
$<_{\text{revlex}}$	reverse lexicographical ordering . . . . .	134

$\text{coef}(X^\alpha, P)$	coefficient of the monomial $X^\alpha$ in the polynomial $P$ . . .	134
$\text{lmon}(P)$	leading monomial of $P$ . . . . .	134
$\text{lcoef}(P)$	leading coefficient of $P$ . . . . .	134
$\text{lt}(P)$	leading term of $P$ . . . . .	134
$\text{Red}(P, X^\alpha, G)$	reduction of $(P, X^\alpha)$ by $G$ . . . . .	134
$\text{Zer}(I, L^k)$	zero set of the ideal $I$ in $L^k$ . . . . .	136
$\text{Proj}_{X_k}(\mathcal{P})$	finite set of polynomials defining the the projection of $\text{Zer}(\mathcal{P}, \mathbb{C}^k)$ to $\mathbb{C}^{k-1}$ . . . . .	137
$A$	$\mathbb{K}[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, \mathbb{K})$ , in Section 4.5 . . . . .	143
$\bar{A}$	$\mathbb{C}[X_1, \dots, X_k]/\text{Ideal}(\mathcal{P}, \mathbb{C})$ , in Section 4.5 . . . . .	144
$S^{-1}A$	ring of fractions . . . . .	146
$\bar{A}_x$	localization of $\bar{A}$ at $x$ . . . . .	146
$e_x$	idempotent associated to $x$ . . . . .	147
$L_f$	map of multiplication by $f$ in the multivariate case . .	149
$\chi(\mathcal{P}, f, T)$	characteristic polynomial of $L_f$ . . . . .	150
$\text{her}(\mathcal{P}, Q)$	Hermite's bilinear map in the multivariate case . . . .	150
$\text{Her}(\mathcal{P}, Q)$	Hermite's quadratic form in the multivariate case . . .	150
$\text{Rad}(\Phi)$	radical of a quadratic form $\Phi$ . . . . .	151
$\text{TaQ}(Q, \mathcal{P})$	Tarski-query of $Q$ for $\mathcal{P}$ . . . . .	151
$\mathbb{P}_k(\mathbb{C})$	complex projective space of dimension $k$ . . . . .	153
$\bar{x} = (x_0: x_1: \dots: x_k)$	homogeneous coordinates of $\bar{x}$ . . . . .	153
$\text{Zer}(\mathcal{P}, \prod_{i=1}^m \mathbb{P}_{k_i}(\mathbb{C}))$	algebraic set of $\mathbb{P}_{k_1}(\mathbb{C}) \times \dots \times \mathbb{P}_{k_m}(\mathbb{C})$ . . . . .	153

## Chapter 5

$D(z, r)$	open disk centered at $z$ with radius $r$ . . . . .	163
$\text{Elim}_{X_k}(\mathcal{P})$	finite set of polynomials obtained by eliminating $X_k$ from a family of polynomials $\mathcal{P}$ for use in cylindrical algebraic decomposition . . . . .	166
$\text{dim}(S)$	dimension of a semi-algebraic set $S$ . . . . .	170
$\bar{\sigma}$	relaxation of the sign condition $\sigma$ . . . . .	173
$\text{Real}(\tau)$	realization of the weak sign condition $\tau$ . . . . .	173
$[a_0, \dots, a_d]$	$d$ -simplex with vertices $a_0, \dots, a_d$ . . . . .	181
$s' \prec s$	$s'$ is face of $s$ . . . . .	181
$s^\circ$	open simplex associated to a simplex $s$ . . . . .	181
$\text{ba}(s)$	barycenter of the simplex $s$ . . . . .	182
$ K $	realization of a simplicial complex $K$ . . . . .	182
$\text{ba}(K)$	barycentric subdivision of a simplicial complex $K$ . . .	182
$\text{dim}(S_x)$	local dimension of $S$ at a point $x \in S$ . . . . .	191

## Chapter 6

$C_p(K)$	$p$ -chain group of a simplicial complex $K$ . . . . .	196
$\partial_p(s)$	boundary of a $p$ -simplex $s$ . . . . .	196
$C_\bullet(K)$	chain complex of a simplicial complex $K$ . . . . .	196

$B_p(\mathbf{C}_\bullet)$	$p$ -boundaries of a chain complex $\mathbf{C}_\bullet$ . . . . .	197
$Z_p(\mathbf{C}_\bullet)$	$p$ -cycles of a chain complex $\mathbf{C}_\bullet$ . . . . .	197
$H_p(\mathbf{C}_\bullet)$	$p$ -th homology group of a chain complex $\mathbf{C}_\bullet$ . . . . .	197
$H_p(K)$	$p$ -th homology group of $\mathbf{C}_\bullet(K)$ . . . . .	198
$H_*(K)$	homology of the chain complex $\mathbf{C}_\bullet(K)$ . . . . .	198
$b_p(K)$	$p$ -th Betti number of a simplicial complex $K$ . . . . .	198
$\chi(K)$	Euler-Poincaré characteristic of $K$ . . . . .	198
$C^p(K)$	$p$ -cochain group of a simplicial complex $K$ . . . . .	199
$\mathbf{C}^\bullet(K)$	cochain complex of a simplicial complex $K$ . . . . .	199
$H^p(\mathbf{C}^\bullet)$	$p$ -th cohomology group of a cochain complex $\mathbf{C}^\bullet$ . . . . .	200
$H^*(\mathbf{C}^\bullet)$	cohomology of a cochain complex $\mathbf{C}^\bullet$ . . . . .	200
$\phi \sim \psi$	$\phi$ and $\psi$ are chain homotopic . . . . .	209
$K^{(n)}$	$n$ -th iterated barycentric subdivision of $K$ . . . . .	213
$\text{star}(a) \subset  K $	union of the relative interiors of all simplices in $K$ having $a$ as a vertex . . . . .	215
$K < L$	$K$ is finer than $L$ . . . . .	215
$\text{diam}(S)$	diameter of $S$ . . . . .	215
$\text{mesh}(K)$	$\max \{\text{diam}(s)   s \in K\}$ for $K$ a simplicial complex . . . . .	216
$H_p(S)$	$p$ -th homology group of a closed and bounded semi-algebraic set $S$ . . . . .	221
$H_*(S)$	homology of a closed and bounded semi-algebraic set $S$ . . . . .	221
$b_p(S)$	$p$ -th Betti number of a closed and bounded semi-algebraic set $S$ . . . . .	221
$b(S)$	sum of the Betti numbers of a closed and bounded semi-algebraic set $S$ . . . . .	221
$\chi(S)$	Euler-Poincaré characteristic of a closed and bounded semi-algebraic set $S$ . . . . .	221
$f \sim g$	$f, g: X \rightarrow Y$ are homotopic . . . . .	223
$f \sim_{sa} g$	$f, g: X \rightarrow Y$ are semi-algebraically homotopic . . . . .	225
$H_*(\text{Reali}(\sigma))$	homology of a sign condition $\sigma$ . . . . .	228
$C_p(K, A)$	$p$ -th chain group of the pair $(K, A)$ . . . . .	228
$H_p(K, A)$	$p$ -th simplicial homology group of the pair $(K, A)$ . . . . .	229
$b_p(K, A)$	$p$ -th Betti number of the pair $(K, A)$ . . . . .	227
$\chi(K, A)$	Euler-Poincaré characteristic of the pair $(K, A)$ . . . . .	229
$H_p(S, T)$	$p$ -th simplicial homology group of a pair of closed and bounded semi-algebraic sets $(S, T)$ . . . . .	229
$b_p(S, T)$	$p$ -th Betti number of a pair of closed and bounded semi-algebraic sets $(S, T)$ . . . . .	229
$\chi(S, T)$	Euler-Poincaré characteristic of a pair of closed and bounded semi-algebraic sets $(S, T)$ . . . . .	229
$H_p^{BM}(S)$	$p$ -th Borel-Moore homology group of a locally closed semi-algebraic set $S$ . . . . .	231

$H_*^{BM}(S)$	Borel-Moore homology of a locally closed semi-algebraic set $S$ . . . . .	231
$b_p^{BM}(S)$	$p$ -th Borel-Moore Betti number of of a locally closed semi-algebraic set $S$ . . . . .	231
$\chi(S)$	Euler-Poincaré characteristic of a locally closed semi-algebraic set $S$ . . . . .	234
$\dot{X}$	Alexandroff compactification of $X$ . . . . .	232
$\text{Reali}(P=0, S)$	$\{x \in S \mid P(x) = 0\}$ . . . . .	235
$\text{Reali}(P > 0, S)$	$\{x \in S \mid P(x) > 0\}$ . . . . .	235
$\text{Reali}(P < 0, S)$	$\{x \in S \mid P(x) < 0\}$ . . . . .	235
$\text{EuQ}(P, S)$	Euler-Poincaré-query of $P$ for $S$ . . . . .	235

**Chapter 7**

$\text{Grad}(Q)(p)$	gradient of $Q$ at $p \in \text{Zer}(Q, \mathbb{R}^k)$ . . . . .	237
$\pi$	projection sending $(x_1, \dots, x_k)$ to $x_1$ . . . . .	238
$S_X$	$S \cap \pi^{-1}(X)$ . . . . .	238
$S_x$	$S_{\{x\}}$ . . . . .	238
$S_{<x}$	$S_{(-\infty, x)}$ . . . . .	238
$S_{\leq x}$	$S_{(-\infty, x]}$ . . . . .	238
$b(k, d)$	maximum of the sum of the Betti numbers of an algebraic set defined by polynomials of degree $d$ in $\mathbb{R}^k$ . . . . .	256
$\text{SIGN}(\mathcal{P})$	set of all realizable sign conditions for $\mathcal{P}$ over $\mathbb{R}^k$ . . . . .	262
$\text{SIGN}(\mathcal{P}, \mathcal{Q})$	set of all realizable sign conditions for $\mathcal{P}$ over $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ . . . . .	262
$\text{Reali}(\sigma, Z)$	$\{x \in \mathbb{R}^k \mid \bigwedge_{Q \in \mathcal{Q}} Q(x) = 0 \wedge \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma(P)\}$ , for $\sigma \in \text{SIGN}(\mathcal{P}, \mathcal{Q})$ . . . . .	262
$b_i(\sigma)$	$i$ -th Betti number of $\text{Reali}(\sigma, Z)$ . . . . .	262
$b_i(\mathcal{Q}, \mathcal{P})$	$\sum_{\sigma} b_i(\sigma)$ . . . . .	262
$\text{deg}(\mathcal{Q})$	maximum of the degrees of the elements of $\mathcal{Q}$ . . . . .	262
$b_i(d, k, k', s)$	maximum of $b_i(\mathcal{Q}, \mathcal{P})$ over all $\mathcal{Q}, \mathcal{P}$ , $\text{deg}(\mathcal{Q}, \mathcal{P}) \leq d$ , $\#(\mathcal{P}) = s$ , and $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ has dimension $k'$ . . . . .	262
$\text{Reali}(\Phi)$	realization of a $(\mathcal{Q}, \mathcal{P})$ -closed formula $\Phi$ . . . . .	268
$b(\Phi)$	sum of the Betti numbers of $\text{Reali}(\Phi)$ . . . . .	268
$b(d, k, k', s)$	maximum of $b(\Phi)$ , over all $(\mathcal{Q}, \mathcal{P})$ -closed formula $\Phi$ , $\text{deg}(\mathcal{Q}, \mathcal{P}) \leq d$ , $\#(\mathcal{P}) = s$ , and $\text{Zer}(\mathcal{Q}, \mathbb{R}^k)$ has dimension $k'$ . . . . .	268
$\text{SIGN}(S)$	set of realizable sign conditions of $\mathcal{P}$ whose realizations are contained in $S$ . . . . .	273
$b(S)$	sum of the Betti numbers of a semi-algebraic set $S$ . . . . .	279

**Chapter 8**

$f(v)$ is $h(O(g(v)))$	there exists a natural number $b$ such that for all $v \in \mathbb{N}^\ell$ , $f(v) \leq h(b(g(v)))$ . . . . .	284
------------------------	--	-----

$f(v)$ is $h(\tilde{O}(g(v)))$	there exist a natural number $a$ such that for all $v \in \mathbb{N}^\ell$ , $f(v) \leq h(bg(v) \log_2(g(v))^a)$ . . . . .	284
$\text{Hor}_i(P, X)$	$i$ -th Horner polynomial associated to $P$ . . . . .	287
$\text{CharPol}(M)$	characteristic polynomial of $M$ . . . . .	299
$\mathcal{F}_n$	polynomials whose degrees are less than $n$ . . . . .	299
$\mathcal{B}$	$\mathcal{B} = X^{n-1}, \dots, X, 1$ , basis of $\mathcal{F}_n$ . . . . .	303
$\text{Mat}(\mathcal{P})$	matrix whose rows are the coordinates of the polynomials in $\mathcal{P}$ in the basis $\mathcal{B}$ . . . . .	303
$\text{pdet}_{m,n}$	$(m, n)$ -polynomial determinant mapping . . . . .	304
$\text{Mat}(\mathcal{P})^*$	$m \times m$ matrix whose first $m - 1$ columns are columns of $\text{Mat}(\mathcal{P})$ and such that the elements of the last column are the polynomials in $\mathcal{P}$ . . . . .	305
$\text{sResP}_j(P, Q)$	$j$ -th signed subresultant polynomial of the polynomials $P$ and $Q$ . . . . .	306
$\text{sResU}_j(P, Q)$	$j$ -th subresultant cofactor of $P$ . . . . .	312
$\text{sResV}_j(P, Q)$	$j$ -th subresultant cofactor of $Q$ . . . . .	312
$\text{sResP}(P, Q)$	signed subresultant sequence of $P$ and $Q$ . . . . .	318

**Chapter 9**

$\text{Bez}(P, Q)$	Bezoutian of $P$ and $Q$ . . . . .	326
$\text{MVar}(s)$	modified number of sign variations in $s$ . . . . .	330
$\text{Han}(\bar{s}_n)$	Hankel matrix of dimension $n \times n$ associated to an infinite sequence, $\bar{s} = s_0, \dots, s_n, \dots$ . . . . .	334
$\text{han}(\bar{s}_n)$	determinant of $\text{Han}(\bar{s}_n)$ . . . . .	334
$\ell_P$	Kronecker form defined on $\mathbb{K}[X]/(P(X))$ . . . . .	335
$\text{Han}(\bar{s}_n)$	Hankel quadratic form associated to $\text{Han}(\bar{s}_n)$ . . . . .	337
$\text{To}(\bar{v})$	triangular Toeplitz matrix associated to $\bar{v}$ . . . . .	339

**Chapter 10**

$C(P)$	$\sum_{q \leq i \leq p} \left  \frac{a_i}{a_p} \right $ , for $P = a_p X^p + \dots + a_q X^q, p > q$ . . . . .	351
$c(P)$	$\left( \sum_{q \leq i \leq p} \left  \frac{a_i}{a_q} \right  \right)^{-1}$ . . . . .	351
$C'(P)$	$(p + 1) \cdot \sum_{q \leq i \leq p} \frac{a_i^2}{a_p^2}$ . . . . .	352
$c'(P)$	$\left( (p + 1) \cdot \sum_{q \leq i \leq p} \frac{a_i^2}{a_p^2} \right)^{-1}$ . . . . .	352
$\ P\ $	norm of a polynomial $P$ . . . . .	353
$\text{Len}(P)$	length of a polynomial $P$ . . . . .	353
$\text{Mea}(P)$	measure of a polynomial $P$ . . . . .	353
$\text{cont}(P)$	greatest common divisor of the coefficients of $P$ , for $P \in$ $\mathbb{Z}[X]$ . . . . .	356
$\text{sep}(P)$	minimal distance between the roots of $P$ . . . . .	
$\text{Bern}_{p,i}(\ell, r)$	$\binom{p}{i} \frac{(X - \ell)^i (r - X)^{p-i}}{(r - \ell)^p}$ , Bernstein polynomials of degree $p$ for $\ell, r$ . . . . .	361

$\text{Rec}_p(P(X))$	reciprocal polynomial . . . . .	361
$\text{Co}_\lambda(P(X))$	contraction by ratio $\lambda$ . . . . .	361
$\text{T}_c(P(X))$	translation by $c$ . . . . .	361
$\text{num}(P; (\ell, r))$	number of roots of $P$ in $(\ell, r)$ . . . . .	362
$\tilde{a}$	list obtained by reversing the list $a$ . . . . .	366
$b(P, \ell, r)$	list of coefficients of $P$ in the Bernstein basis of $\ell, r$ , . .	371
$\mathcal{C}(\ell, r)$	closed disk with diameter $(\ell, r)$ . . . . .	371
$\mathcal{C}_1(\ell, r)$	closed disk defined by the circle circumscribed to the equilateral triangle $T_1$ based on $[\ell, r]$ (see Figure 10.5) . . .	371
$\mathcal{C}_2(\ell, r)$	closed disk symmetric to $\mathcal{C}_1$ with respect the $X$ -axis . .	371
$\tilde{P}$	separable part of the polynomial $P$ . . . . .	372
$P[\ell, r]$	$2^{pk}\text{Co}_{r-\ell}(\text{T}_{-\ell}(P))$ . . . . .	378
$w(\ell, r)$	$r - \ell$ , width of the interval $(\ell, r)$ . . . . .	375
$c(P = 0, Z)$	number of elements in $\text{Reali}(P = 0, Z)$ . . . . .	383
$c(P > 0, Z)$	number of elements in $\text{Reali}(P > 0, Z)$ . . . . .	383
$c(P < 0, Z)$	number of elements in $\text{Reali}(P < 0, Z)$ . . . . .	383
$\text{TaQ}(P, Z)$	Tarski-query of $P$ for $Z$ . . . . .	383
$\text{Reali}(\sigma, Z)$	$\{x \in Z \mid \bigwedge_{P \in \mathcal{P}} \text{sign}(P(x)) = \sigma(P)\}$ . . . . .	383
$c(\sigma, Z)$	$\#(\text{Reali}(\sigma, Z))$ . . . . .	383
$\text{SIGN}(\mathcal{P}, Z)$	list of sign conditions realized by $\mathcal{P}$ on $Z$ . . . . .	384
$c(\mathcal{P}, Z)$	list of $c(\sigma, Z)$ for $\sigma \in \text{SIGN}(\mathcal{P}, Z)$ . . . . .	384
$\text{Mat}(A, \Sigma)$	matrix of signs of $\mathcal{Q}^A$ on $\Sigma$ . . . . .	384
$\text{Ada}(\mathcal{P}, Z)$	list of polynomials adapted to sign determination . . .	387

**Chapter 11**

$\mathcal{C}_i(\mathcal{P})$	$i$ -th level cylindrifying family of polynomials associated to a family of polynomials $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ . . . . .	404
$\mathcal{C}(\mathcal{P})$	cylindrifying family of polynomials associated to $\mathcal{P}$ . .	404
$\text{CSIGN}(\mathcal{P})$	tree of cylindrical realizable sign conditions of a family of polynomials $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ . . . . .	416
$\text{CSIGN}(\mathcal{P}, \mathcal{A})$	tree of cylindrical realizable sign conditions of $\mathcal{P}$ on $\mathcal{A}$ . . . . .	417
$\text{size}(\Phi)$	size of a formula $\Phi$ . . . . .	427
$\overline{\mathcal{C}_i(\mathcal{P})}$	result of saturating $\text{Elim}_{X_{i+1}}(\overline{\mathcal{C}_{i+1}(\mathcal{P})})$ by derivation with respect to $X_i$ . . . . .	425
$\text{RElim}_{X_k}(P, \mathcal{P})$	result of Restricted Elimination, eliminating $X_k$ from a polynomial $P$ and a family of polynomials $\mathcal{P}$ . . . . .	441

**Chapter 12**

$\text{S}(P_1, P_2)$	S-Polynomial of $P_1$ and $P_2$ . . . . .	446
$\text{Mon}(\mathcal{G})$	set of monomials under the staircase . . . . .	449
$\text{Cor}(\mathcal{G})$	corners of the staircase of $\mathcal{G}$ . . . . .	449
$\text{Bor}(\mathcal{G})$	border of the staircase of $\mathcal{G}$ . . . . .	449

NF( $P$ )	normal form	450
Tab( $\mathcal{B}$ )	$\{a b, b \text{ in } \mathcal{B}\}$	451
Mat( $\mathcal{B}$ )	multiplication table of $A$ in $\mathcal{B}$	451
$\chi(a, T)$	characteristic polynomial of $L_a$	463
$\varphi(a, f, T)$	$\sum_{x \in \text{Zer}(\mathcal{P}, \mathcal{C}^k)} \mu(x) f(x) \prod_{t \in \text{Zer}(\chi(a, T), \mathcal{C})} (T - t)$	463
$x_u(t)$	point associated to a univariate representation $u$	465
$\varphi(a, T)$	$(\varphi(a, 1, T), \varphi(a, X_1, T), \dots, \varphi(a, X_k, T))$	465
$\varphi_b(a, T)$	$(\varphi(a, b, T), \varphi(a, b X_1, T), \dots, \varphi(a, b X_k, T))$	466
$<_\varepsilon$	order on $\mathbb{Q}^m$ defined by $\varepsilon^\nu <_\varepsilon \varepsilon^\mu$ if and only if $(\nu_m, \dots, \nu_1) >_{\text{lex}} (\mu_m, \dots, \mu_1)$	471
$K(\varepsilon)_b$	elements of $K(\varepsilon)$ which are sums of $\varepsilon^\nu$ with $\varepsilon^\nu \leq_\varepsilon 1$	471
$K\langle \varepsilon \rangle_b$	elements of $K\langle \varepsilon \rangle$ which are sums of $\varepsilon^\nu$ with $\varepsilon^\nu \leq_\varepsilon 1$	471
$o(\tau)$	order of $\tau$	471
In( $\tau$ )	initial coefficient of $\tau$	471
$\lim_\varepsilon(\tau)$	limit map from $K\langle \varepsilon \rangle_b$ to $K$	471
$\text{Zer}_b(F(T), C\langle \varepsilon \rangle)$	$\{\tau \in C\langle \varepsilon \rangle_b \mid F(\tau) = 0\}$	472
$\text{Zer}_b(\mathcal{P}, R\langle \varepsilon \rangle^k)$	$\text{Zer}(\mathcal{P}, R\langle \varepsilon \rangle^k) \cap R\langle \varepsilon \rangle_b^k$	473
$\text{Zer}_b(\mathcal{P}, C\langle \varepsilon \rangle^k)$	$\text{Zer}(\mathcal{P}, C\langle \varepsilon \rangle^k) \cap C\langle \varepsilon \rangle_b^k$	473
tDeg $_{X_i}(Q)$	total degree of $Q$ in $X_i$	485
$G_k(\vec{d}, c)$	$c^{\vec{d}_1}(X_1^{\vec{d}_1} + \dots + X_k^{\vec{d}_k} + X_2^2 + \dots + X_k^2) - (2k - 1)$	485
Def( $Q, \zeta$ )	$\zeta G_k(\vec{d}, c) + (1 - \zeta)Q$	485
Cr( $Q, \zeta$ )	$\{\text{Def}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k}\}$	487
Def $_+(Q, \zeta)$	$\text{Def}(Q, \zeta) + X_{k+1}^2$	487
Cr $_+(Q, \zeta)$	$\{\text{Def}(Q, \zeta), \frac{\partial \text{Def}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \text{Def}(Q, \zeta)}{\partial X_k}, 2X_{k+1}\}$	487
$\overline{\text{Def}}(Q, \zeta)$	reduction of $\text{Def}(Q, \zeta)$	491
$\overline{\text{Cr}}(Q, \zeta)$	$\{\overline{\text{Def}}(Q, \zeta), \frac{\partial \overline{\text{Def}}(Q, \zeta)}{\partial X_2}, \dots, \frac{\partial \overline{\text{Def}}(Q, \zeta)}{\partial X_k}\}$	492
$Q_u$	$g_0^e Q\left(\frac{g_k}{g_0}, \dots, \frac{g_k}{g_0}\right)$ for $Q \in K[X_1, \dots, X_k]$ and $u = (f, g_0, \dots, g_k)$ a $k$ -univariate representation	495

**Chapter 13**

$H_k(d, i)$	polynomial $1 + \sum_{1 \leq j \leq k} i^j X_j^d$	508
$H_k^h(d, i)$	homogeneous polynomial $X_0^d + \sum_{1 \leq j \leq k} i^j X_j^d$	508
$P_u$	$g_0^e P\left(\frac{g_k}{g_0}, \dots, \frac{g_k}{g_0}\right)$ for $P \in K[X_1, \dots, X_k]$ and $u = (f, g_0, \dots, g_k)$ a $k$ -univariate representation	519
$v_k(x)$	Vandermonde vector $(1, x, \dots, x^{k-1})$	521
$V_\ell$	vector space generated by $v_k(\ell), v_k(\ell + 1), \dots, v_k(\ell + k - k' - 1)$	521
$\mathcal{L}_{k, k-k'}$	set of of $k - k'$ -planes such that any plane of dimension $k'$ is transversal to at least one element of the family $\mathcal{L}_{k, k-k'}$	521

**Chapter 14**

$\text{SIGN}_{\Pi}(\mathcal{P})$	tree of realizable sign conditions of $\mathcal{P}$ with respect to a partition $\Pi$ of the variables . . . . .	534
$\text{SIGN}_{\Pi}(\mathcal{P}, \mathcal{A})$	tree of realizable sign conditions of $\mathcal{P}$ for $\mathcal{A}$ with respect to $\Pi$ , where $\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_\omega$ is a $\Pi$ -set . . . . .	537
$\text{Remo}_{\varepsilon}(\Phi)$	result of removing $\varepsilon$ from the formula $\Phi$ . . . . .	541
$\text{B}_{\Pi, i}(\mathcal{P})$	result after eliminating all but the first $i$ blocks in the Block Elimination algorithm . . . . .	542
$\text{SSIGN}(\mathcal{P})$	set of strict realizable sign conditions of $\mathcal{P}$ i.e. the realizable sign conditions $\sigma \in \{0, 1, -1\}^{\mathcal{P}}$ such that for every $P \in \mathcal{P}$ , $P \neq 0$ , $\sigma(P) \neq 0$ . . . . .	558

**Chapter 15**

$\text{RM}(\text{Zer}(Q, \mathbb{R}^k), \mathcal{N})$	roadmap for the algebraic set $\text{Zer}(Q, \mathbb{R}^k)$ passing through the finite set of points associated to $\mathcal{N}$ . . . . .	576
---	--	-----

**Chapter 16**

$\text{URM}(Q, \mathcal{P})$	uniform roadmap for a family $\mathcal{P} \subset \mathbb{R}[X_1, \dots, X_k]$ of polynomials contained in the algebraic set $\text{Zer}(Q, \mathbb{R}^k)$ . . .	604
$\text{RM}(S)$	roadmap for the semi-algebraic set $S$ . . . . .	605



---

# Index

- Absolute value, 32
- Active interval, 371
- Acyclic, 211
- Addition
  - Algorithm
    - of matrices, 290
    - of polynomials, 282, 285
- Additivity of Euler-Poincaré characteristic, 235
- Affinely independent, 181
- Alexandroff compactification, 232
- Algebraic
  - closure, 11
  - set, 11, 57, 136, 153
- Algorithm, 281
- Alternating mapping, 303
- Approximating varieties, 523
- Archimedean, 43, 43
- Ascending sequence of simplices, 182
- Atom, 12, 58
  - $\mathcal{P}$  -, 415
- Bézout, 157
- Ball
  - closed -, 83
  - open -, 83
- Band, 160
- Barycenter, 182
- Barycentric subdivision, 182
- Bernstein coefficients
  - Algorithm, 365
  - Special, 367
- Bernstein polynomials, 361
- Betti numbers
  - of a closed and bounded semi-algebraic set, 221
  - of a pair, 229, 229
  - of a simplex, 198
  - of algebraic sets, 257
  - of closed semi-algebraic sets, 268
  - of semi-algebraic sets, 279
  - of sign conditions, 262
- Bezoutian, 326
- Bitsize, 285
- Block elimination
  - Algorithm, 539
- Block structure, 534
- Block structured signs
  - Algorithm, 543
- Boundary, 196, 197
- Buchberger
  - Algorithm, 447
- Budan-Fourier theorem, 45
- Carrier function, 212
  - acyclic, 212
- Cauchy index, 53
  - Algorithm, 325
  - Sturm, 323
- Cauchy-Binet formula, 104
- Cauchy's bound, 352
- Cauchy-Schwarz inequality, 123
- Cell, 159
- Chain, 196
  - complex, 197
  - group, 196
  - homotopy, 209

- map, 197
- Characteristic of a field, 26
- Characteristic polynomial
  - Algorithm, 299
- Closed, 84
  - in  $S$ , 84
- Closed under differentiation, 173
- Closure, 84
  - of a semi-algebraic set, 84
- Co-boundary, 200
- Cochain
  - complex, 199, 199
- Co-chain
  - map, 199
- Co-cycle, 200
- Coefficient, 15
- Cohomology, 200
- Comparison of roots
  - Algorithm
    - Parametrized, 581
    - Real, 382
    - Real Closed Field, 399
    - Real Triangular, 411
    - Recursive, 413
    - Triangular, 496
- Complexity, 282
  - binary, 285
- Computing realizable sign conditions
  - Algorithm, 511
- Cone, 33
  - positive -, 33
  - proper -, 33
- Congruence, 132
- Conic structure
  - at infinity, 189
  - local, 189
- Conjugate, 40
- Connected, 86
  - semi-algebraically -, 86
  - semi-algebraically path -, 86, 156
- Connected components
  - Algorithm
    - Algebraic Set, 591
    - Basic, 612
    - Semi-algebraic Set, 616
    - semi-algebraically -, 168
- Connecting
  - Algorithm
    - Algebraic, 579
- Bounded, 601
- Bounded Algebraic, 577
- Parametrized Bounded, 608
- Parametrized Bounded Algebraic, 590
- Constant rank, 192
- Constructible
  - basic - set, 12
  - set, 12, 57
- Content, 356
- Contiguous, 212
- Continuity, 84
- Continuity of roots, 164
- Continuous extension, 92
- Control line, 363
- Control polygon, 363
- Convex, 87
- Coprime, 15
- Covering by Contractible Sets
  - Algorithm, 630, 631
- Critical
  - ordinary - point, 432
  - point, 191, 238, 432
    - non-degenerate, 243
  - value, 191, 238
- Critical point method, 237, 445, 505, 533, 563
- Curve segment, 431, 573
  - Algorithm, 573
  - Modified, 618
  - Parametrized, 584
  - representation, 572
  - parametrized, 583
- Curve selection lemma, 92
- Cycle, 197
- Cylindrical decomposition, 159
  - adapted to  $\mathcal{P}$ , 163
  - Algorithm, 406
  - Improved, 425
  - Stratifying, 428
  - induced, 162
- Data for adjacencies, 614
  - Algorithm, 614
- Decision
  - Algorithm
    - Cylindrical, 422
    - General, 546
  - problem, 403

- Deformation retraction, 189, 224
  - semi-algebraic -, 225
- Degree, 15, 136, 485
- Derivative, 29, 95, 99
- Descartes law of signs, 44
- Descartes' Real root isolation
  - Algorithm, 379
- Diameter, 215
- Dickson lemma, 133
- Diffeomorphism
  - $S^\infty$  -, 97
- Differentiable, 94
- Dimension
  - Algorithm, 559
  - local, 191
  - of a  $S^\infty$  submanifold of  $\mathbb{R}^k$ , 97
  - of a semi-algebraic set, 170
  - of a simplex, 181
  - of complex projective space, 153
- Discriminant, 40, 101, 102
- Distinguished
  - hyperplane, 569
  - point, 569
  - value, 569
- Divergence property, 578
- Divisor, 15
  - greatest common -, 15
  - greatest common , 19
- Dogson-Jordan-Bareiss
  - Algorithm, 297
- Edge, 181
- Elimination
  - Algorithm, 404
  - Restricted, 441
- Euclidean division, 15
  - Algorithm, 284
- Euclidean norm, 83
- Euler-Poincaré characteristic
  - Algorithm
    - Algebraic Set, 502
    - Bounded Algebraic Set, 499
    - Sign Conditions, 531
  - of a closed and bounded semi-algebraic set, 221
  - of a locally closed semi-algebraic set, 234
  - of a pair, 229, 229
  - of a simplicial simplex, 198
- Euler-Poincaré-query, 235
  - Algorithm, 530
- Evaluation of polynomials
  - Algorithm, 288, 288
- Exact division, 282
  - Algorithm
    - Multivariate polynomials, 287
- Exact sequence, 206
  - short, 206
- Existential theory of the reals, 505, 516
  - Hardness, 513
- Extension, 27, 72, 73
- Face, 181
  - proper-, 181
- Facet, 181
- Factorization, 40
- Field
  - algebraically closed -, 11
  - ordered -, 31, 84
  - real, 33
  - real closed -, 34
- Finite mapping, 137, 139
- First Betti number
  - Algorithm, 632
- Flat, 97
- Flow line, 239
- Formula
  - $\mathcal{P}$  -, 415
  - basic-, 13, 59
  - $(\mathcal{Q}, \mathcal{P})$ -closed, 268
  - of the language of fields, 12
  - of the language of ordered fields, 58
  - quantifier free -, 13, 59
- Fundamental theorem of algebra, 39
- Gauss
  - Algorithm, 294
- Gcd and Gcd-free part
  - Algorithm, 356
- Gcd-free part, 355
- General position
  - $\ell$  -, 508, 508
  - strong  $\ell$  -, 508
- Generalized Permanences minus Variations
  - Algorithm, 324
- Generic position, 432
- Germes of semi-algebraic continuous functions, 88
- Global Optimization
  - Algorithm, 557

- Gröbner basis, 135, 135
  - special, 456
  - special parametrized -, 458
- Gradient, 237
  - flow, 240
- Gram-Schmidt orthogonalization, 122
- Hadamard bound, 292
- Hankel matrix, 334
- Hankel quadratic forms, 337
- Hardt triviality, 186
- Hermite
  - bilinear map, 150
  - quadratic form, 128, 150
- Hermite theorem, 130
  - multivariate -, 151
- Hessian, 243
- Hilbert basis theorem, 132
- Hilbert Nullstellensatz, 136, 140, 140, 142
- Homeomorphism
  - semi-algebraic -, 86
- Homogeneous coordinates, 153
- Homogeneous polynomial, 136
- Homology, 197, 198, 221, 226
  - Borel-Moore, 231
  - Borel-Moore - groups, 231
  - simplicial - groups, 198, 221
  - of a pair, 229, 229
- Homotopy
  - semi-algebraic -, 225
- Homotopy equivalence, 224
  - of pairs, 230
  - semi-algebraic -, 225
- Horner polynomials, 287, 334
- Hypersurface
  - non-singular algebraic -, 237
- Ideal, 132
  - generated by a set of polynomials, 132
  - prime -, 132
  - principal -, 132
- Idempotent, 147
- Implicit function theorem, 97
  - projective -, 156
- Index, 243
- Infinitesimal, 32, 75
- Inner product, 122
- Interior, 84
- Intermediate points
  - Algorithm, 400
  - Parametrized, 581
  - Recursive, 414
  - Triangular, 497
- Intermediate value, 34, 85
- Invariant
  - $\mathcal{P}$ -, 163
- Inverse function theorem, 96
- Inverse sign determination
  - Algorithm, 547
- Isolating list, 373, 380
- Isolating parallelepiped, 409
- Isolating tree, 374
- Jacobian, 95
  - matrix, 95
- Jump of a function, 53
- Kronecker form, 335
- Language
  - of fields, 12, 12
  - of orderd fields, 58
  - of ordered fields, 58
- Leading
  - coefficient, 15, 134
  - monomial, 134
  - term, 134
- Leading vertex, 182
- Lefschetz principle, 26
- Length, 353
- Lexicographical ordering, 35
  - graded, 36
  - reverse, 134
- Lifting phase
  - Algorithm, 415
  - Real, 412
- Limit
  - Algorithm
    - Parametrized Real Bounded Points, 481
    - Real Bounded Points, 480
    - of a bounded Puiseux series, 81
- Limit of a Puiseux series, 81
- Linear recurrent sequence, 334
- Linking paths
  - Algorithm, 602
- Linking points
  - Algorithm, 621
- Local ring, 146
- Localization, 146
- Matrix of signs, 64, 384

- Mayer-Vietoris, 209, 222, 227
- Mean value theorem, 41
- Measure, 353
- Mesh, 216
- Minimal distance between roots, 357
- Minimum of a polynomial, 551
- Modified sign variations, 330
- Modulus, 40
- Monomial, 134
  - Ordering, 133
- Morse
  - function, 243
  - lemma A, 239
  - lemma B, 247
- Multihomogeneous polynomial, 153
- Multilinear mapping, 303
- Multiple
  - least common -, 15
- Multiplication
  - Algorithm
    - of matrices, 291, 292
    - of polynomials, 283, 286
- Multiplication map, 128, 149, 454
- Multiplication table, 451
  - Algorithm, 452
  - Parametrized Special, 458
  - Special, 457
  - entries, 452
  - size, 452
- Multiplicity
  - of a root, 30
  - of a zero, 148
  - virtual -, 50
- Nash functions, 95
- Newton
  - diagram, 76
  - polygon, 76
- Newton sums, 102
  - Algorithm, 290
- Norm, 83, 95, 122, 353
- Normal form, 450
  - Algorithm, 450
- Normal polynomial, 47
- NP, 513
  - complete, 513
  - hard, 513
- Number of distinct zeros
  - Algorithm, 454
- Number of sign variations, 44, 330
  - modified -, 331
- Oleinik-Petrovski/Thom/Milnor bound, 257
- Open, 84
  - disk, 163
  - in  $S$ , 84
  - in projective space, 155
- Order, 75, 471
- Ordered set
  - partially, 31
  - totally, 31
- Orthant, 449
- Orthogonal, 122
- Parametrized path, 588
- Partition of a Line
  - Algorithm, 400
- Permanences minus Variations, 113
- Polyhedron, 182
- Polynomial determinant, 304
- Prenex normal form, 14, 59
- Principal minor, 296
- Projection
  - Algorithm, 571
  - Parametrized, 582
- Projection theorem, 25, 60, 68, 154
- Projective
  - complex - space, 153
- Pseudo-critical
  - point, 488
  - value, 488
- PSPACE, 513
- Quadratic form, 119
  - Diagonal expression of a -, 119
- Quantifier elimination, 25, 69
  - Algorithm, 549
  - Cylindrical, 426
  - Local, 554
  - problem, 403
- Quasi-monic, 136
- Quotient, 15
- Quotient ring, 132
- Radical
  - of an ideal, 132
- Rank, 119
- Real algebraic numbers, 34
- Real closure, 43
- Real eigenvalues, 123
- Real root isolation
  - Algorithm, 373

- Triangular, 410
- Realization
  - of a formula, 13, 59
  - of a sign condition, 42, 64, 383
  - of a simplicial complex, 182
  - of a weak sign condition, 528
- Reduction, 134
- Regular
  - point, 191
  - value, 191
- Relaxation
  - of a sign condition, 173
- Remainder, 15
  - extended signed - sequence, 17
- Removal of infinitesimals
  - Algorithm, 443
- Resultant, 106
- Reversing rows, 111
- Ring
  - ordered -, 31
- Ring of fractions, 146
- Roadmap, 563
  - Algorithm
    - Algebraic, 578
    - Bounded Algebraic, 576
    - General, 625
    - Parametrized Bounded Algebraic, 587
    - uniform -, 602
- Rolle theorem, 41
- Root, 11
- Sample points, 511
  - $\Pi$  -, 537, 538
  - Algorithm
    - Real Triangular, 411
    - Real Univariate, 382
    - Recursive, 414
    - Triangular, 498
    - Univariate, 401
  - of a cylindrical decomposition, 406
  - ordered list of -, 401
- Sampling
  - Algorithm, 514
  - Algebraic, 493
  - Bounded Algebraic, 492
  - on an algebraic set, 526
  - Parametrized Bounded Algebraic, 494
- Sard theorem, 192
- Semi-algebraic
  - $\mathcal{P}$  -, 163
  - basic - set, 57
  - basic locally closed, 231
  - continuous function, 85
  - function, 71
  - in projective space, 155
  - set, 57
  - set defined over  $D$ , 57
- Sentence, 14, 59
  - $\mathcal{P}$  -, 415
- Separable, 31, 431
  - part, 355
- Separating, 145
- Series
  - algebraic Puiseux -, 81
  - formal power-, 74
  - Laurent -, 74
  - Puiseux -, 74
- Set of possible gcd's, 24
- Sign, 31
- Sign at a Point
  - Algorithm
    - Triangular Real, 410
- Sign at a Root
  - Algorithm
    - Real, 380
    - Real Closed Field, 398
- Sign condition, 41
  - realizable -, 42, 173
  - refinement, 581
  - set of realizable sign conditions, 262, 525
  - strict -, 41
  - tree of cylindrical realizable , 416
  - tree of realizable - with respect to  $\Pi$ , 534
  - weak -, 173
- Sign determination
  - adapted to -, 386, 530
  - Algorithm, 390
  - Adapted family, 389
  - Family adapted to -, 396
  - Multivariate, 456
  - Naive, 385
  - Parametrized, 551
  - Recursive, 412
  - Triangular, 495

- Univariate, 396
- Signature, 121
- Signature of Hankel form
  - Algorithm, 343
- Signature through Descartes
  - Algorithm, 301
- Signed pseudo-remainder, 21
  - tree of possible - sequences, 22
- Signed remainder sequence
  - Algorithm, 301
  - Extended, 302
- Signed subresultant
  - Algorithm, 318
  - Extended, 320
- Simplex, 181
  - oriented -, 195
- Simplicial
  - approximation, 218
  - complex, 182
  - decomposition, 182
  - map, 196
- Singular point, 432
- Size
  - of signed remainders, 316
  - of signed subresultants, 315
- Size of input, 282
- Smooth point, 97, 191
- Solution
  - set of solutions, 143
- Special values
  - Algorithm, 600
- Sperner map, 214
- Sphere, 83
- S-polynomial, 445
- Square-free, 31
- Stability
  - domain of -, 345
- Staircase
  - border, 449
  - corner, 449
  - monomials under the -, 449
- Star, 215
- Stickelberger, 150
- Stratification
  - cell - adapted to  $\mathcal{P}$ , 177
  - semi-algebraic-, 177
- Stratifying family, 180
- Stratum, 177
- Structure, 281
- Structure theorem for signed subresultants, 307
- Sturm sequence, 52
- Sturm's theorem, 52
- Subdivision, 182
- Submanifold
  - $\mathcal{S}^\infty$  - of  $\mathbb{R}^k$ , 97
- Subresultant
  - defective, 307
  - non-defective, 307
  - signed, 306
  - signed - coefficient, 110
  - signed - sequence, 318
  - signed - transition matrix, 313
- Sums of squares, 33
- Sylvester matrix, 106
- Sylvester's law of inertia, 120
- Sylvester-Habicht matrix, 110
- Symmetric, 35
  - elementary - function, 35
- Tangent space, 97, 97
- Tarski-query, 54, 383
  - Algorithm, 325
  - Multivariate, 455
  - Remainder, 324
- Tarski's theorem, 57
- Tarski-Seidenberg principle, 70, 70
- Taylor formula, 29
- Tensor product, 65
- Term, 134
- Termination
  - criterion, 446
- Theorem of 3 circles, 371
- Thom encoding, 42, 43, 397, 413
  - Algorithm, 397
  - Parametrized Triangular, 553
  - Recursive, 413
  - Triangular, 496
  - Multivariate, 412
  - ordered list of -, 397, 399
  - parametrized triangular, 553, 553
- Thom lemma, 42, 173
  - generalized -, 178
- Toeplitz matrix, 338
- Topological types, 188, 430
- Topology of a curve
  - Algorithm, 436
- Touching points
  - Algorithm, 623

- Trace, 128
  - Algorithm, 454
- Transfer principle, 26, 70, 70
- Translation
  - Algorithm, 289, 289
- Transpose, 119
- Transversal, 520
- Tree of possible signed pseudo-remainder sequences, 22
- Triangular system, 409
  - parametrized, 547
- Triangulation, 183, 183
  - respecting a semi-algebraic family, 183
- Truncation, 21
  - set of -, 22
- Unbounded, 32
- Univariate representation, 465
  - Algorithm
    - Candidate, 468
    - Simple, 470
  - parametrized, 481
  - parametrized real, 582
  - point associated to a real -, 465
  - points associated to a -, 465
  - real -, 465
  - triangular real -, 571
- Valuation ring, 81
- Value
  - special -, 594, 594
- Vandermonde
  - determinant, 104
  - matrix, 104
- Variable
  - bound -, 14, 59
  - free -, 12, 58
- Vector field, 239
- Vertex, 181
- Virtual
  - multiplicity, 50
  - roots, 50
- Well-separating, 473
- Width of an interval, 375
- Zero, 153
  - non-singular -, 148
  - non-singular projective, 155
  - set of -, 11, 57
  - simple, 148
- Zero-dimensional, 143
  - Algorithm
    - Arithmetic Operation, 453
- Zigzag Lemma, 207, 208