



# Theory and Methods of Statistics



# Theory and Methods of Statistics

P. K. Bhattacharya  
Prabir Burman



ELSEVIER

AMSTERDAM • BOSTON • HEIDELBERG • LONDON  
NEW YORK • OXFORD • PARIS • SAN DIEGO  
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO

Academic Press is an imprint of Elsevier



Academic Press is an imprint of Elsevier  
32 Jamestown Road, London NW1 7BY, UK  
525 B Street, Suite 1800, San Diego, CA 92101-4495, USA  
50 Hampshire Street, 5<sup>th</sup> Floor, Cambridge, MA 02139, USA  
The Boulevard, Langford Lane, Kidlington, Oxford OX5 1GB, UK

Copyright © 2016 Elsevier Inc. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Details on how to seek permission, further information about the Publisher's permissions policies and our arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: [www.elsevier.com/permissions](http://www.elsevier.com/permissions).

This book and the individual contributions contained in it are protected under copyright by the Publisher (other than as may be noted herein).

#### **Notices**

Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

#### **British Library Cataloguing in Publication Data**

A catalogue record for this book is available from the British Library

#### **Library of Congress Cataloging-in-Publication Data**

A catalog record for this book is available from the Library of Congress

ISBN: 978-0-12-802440-9

For information on all Academic Press publications  
visit our website at <http://www.elsevier.com/>



Working together  
to grow libraries in  
developing countries

[www.elsevier.com](http://www.elsevier.com) • [www.bookaid.org](http://www.bookaid.org)

*Publisher:* Nikki Levy

*Acquisition Editor:* Graham Nisbet

*Editorial Project Manager:* Susan Ikeda

*Production Project Manager:* Poulouse Joseph

*Designer:* Victoria Pearson

Typeset by SPi Global, India

To my wife, Srilekha Bhattacharya.  
P. K. Bhattacharya

To my mother, Shanti Lata Barman.  
Prabir Burman



# Preface

This book has grown out of lecture notes in courses in Mathematical Statistics, Linear Models, Multivariate Analysis, and Time Series Analysis taught by the authors over many years at UC Davis. Of these, Mathematical Statistics ([Chapters 1–10](#)) and Linear Models ([Chapter 11](#)) are core courses for PhD students in Statistics and Biostatistics at most institutions. Courses in Multivariate Analysis and Time Series Analysis are also taken by many students in these programs. A good background in Advanced Calculus and Matrix Analysis is a prerequisite for these courses.

Although most students in these courses would have taken an intermediate level course in Probability Theory, we have included such material in [Chapters 1 and 2](#) for a review. [Chapter 3](#) introduces various modes of convergence of an infinite sequence of random variables, which may be deferred until the study of asymptotic theory begins in [Chapter 7](#). In [Chapter 4](#), we outline the main problems of statistical inference and various approaches to optimality in the decision theoretic framework before treating Point Estimation, Hypothesis Testing, and Confidence Sets in [Chapters 5 and 6](#). Methods based on Likelihood, Distribution-free Tests, and Curve Estimation are treated in [Chapters 7–9](#), and [Chapter 10](#) deals with Statistical Functionals. Gauss-Markov Models, topics in Model selection and Linear Mixed Models are the main focus of [Chapter 11](#). [Chapter 12](#) deals with Multivariate Analysis covering many of the standard methods used in practice. An introduction to Time Series Analysis is given in [Chapter 13](#) covering aspects of ARMA modeling and Spectral Analysis. These chapters mostly concentrate on asymptotic properties of these methodologies, using the material in [Chapter 3](#). Some technical results are included in [Appendices A and B](#).

Throughout the book we have restricted to discrete and absolutely continuous random variables. Except for using the concept of a  $\sigma$ -field of subsets and in particular the Borel  $\sigma$ -field in Euclidean spaces, we have avoided all technicalities of measure theory. On the other hand, we have used Stieltjes integrals, providing a quick introduction to these integrals in [Appendix A](#). Basic results from matrix algebra and distribution of quadratic forms are included in [Appendix B](#).

We have not been able to cover everything in [Chapters 1 through 10](#) and [Appendix A](#) in a three-quarter course with three lectures and one discussion per week in Mathematical Statistics. The following possibility of slightly trimming this material is suggested so as to fit into a three-quarter framework.

- (i) Give a quick review of [Chapters 1](#) and [2](#), leaving out most of the materials for self-study.
- (ii) Introduce as many topics as possible out of locally best tests ([Section 6.8](#)), SPRT ([Section 6.11](#)), locally most powerful rank tests ([Section 8.2](#)), and curve estimation ([Chapter 9](#)).
- (iii) Leave the details of most of the proofs in [Sections 6.7–6.9](#), [7.2](#), [7.4](#), and [10.6](#) for self-study, going over only the main ideas in class-room presentation.

We also suggest the following outline for a full academic year Master's level course in Probability and Mathematical Statistics.

- (i) [Chapters 1–3](#) omitting [Theorems 3.2.5](#) parts IX–XI; [Theorems 3.2.3–3.2.5](#).
- (ii) Omit the following: proofs in [Sections 4.6.5](#) and [4.6.6](#); Ancillarity and Completeness; proof in [Section 5.2.3](#); Equivariance in its generality; [Section 6.8](#) and [Section 6.11](#); proofs in [Sections 6.12.2](#), [7.1.1](#), and [7.1.2](#); proofs of [Theorems 7.2.1](#), [7.2.2](#), and [7.4.1](#); [Section 8.2](#); proofs in [Section 8.3](#); proofs in [Chapter 9](#) and the entire [Chapter 10](#).

[Chapters 5](#) and [6](#) are influenced by Ferguson [1] and Lehmann [2, 3], and [Chapter 8](#) relies heavily on Hájek and Šidák [4].

The material in [Chapter 11](#) can be fitted in a two-quarter course in Linear Models with minor modifications; introduce as many topics as possible from Model Selection omitting some of the proofs, using examples to clarify the concepts.

Most of the material in [Chapter 12](#) should be accessible to doctoral students along with the examples. An introduction to Time Series ([Chapter 13](#)) can be done by focusing on the main ideas and examples. The material in these two chapters can be fitted in two one-quarter courses for PhD students.

[Chapters 11](#) and [12](#) use the material provided in [Appendix B](#).

Our work on this book was initiated at the insistence of Jerome Braun, who took the Mathematical Statistics course from one of us. Thank you Jerome! The authors are also grateful to Ms Christine Cai for typing some parts of the book. We are indebted to the National Science Foundation for a grant (DMS 09-07622) for support of research in writing this book.

P.K. Bhattacharya, Prabir Burman  
*Davis, California*

# Probability Theory

## 1.1 Random Experiments and Their Outcomes

Probability theory is the systematic study of outcomes of a random experiment such as the roll of a die, or a bridge hand dealt from a thoroughly shuffled deck of cards, or the life of an electric bulb, or the minimum and the maximum temperatures in a city on a certain day, etc. The very first step in such a study is to *visualize* all possible outcomes of the experiment in question, and then to *realize* that the actual outcome of such an experiment is not predictable in advance. However, from the nature of the experiment, or from our experience with the past results of the experiment (if available), we may be able to assign probabilities to the possible outcomes or sets thereof.

For example, in the roll of a balanced, six-faced die, the possible outcomes are  $\{1, 2, 3, 4, 5, 6\}$  to each of which we may assign a probability of  $\frac{1}{6}$  (ie, in many repetitions of the trial, we expect each of these outcomes to occur  $\frac{1}{6}$  of the times). From this, we can also conclude the outcome to be an even number with probability  $\frac{3}{6} = \frac{1}{2}$ .

Similarly, the possible outcomes of a bridge hand dealt from a standard deck are all the  $\binom{52}{13} = \frac{52!}{13!39!}$  combinations of 13 cards from a 52-card deck, each carrying a probability of  $1/\binom{52}{13}$ .

The possible outcomes of the life of an electric bulb (in hours) is the set  $[0, \infty)$  and the possible outcomes of the minimum and maximum temperatures in Philadelphia on a certain day in future (in  $^{\circ}\text{C}$ ) is  $\{(x, y): -273 < x < y < \infty\}$ . In these last two examples, the sets of possible outcomes may seem to be unrealistically large, but this is because we cannot put an upper limit to the life of an electric bulb, or lower and upper limits to the temperatures in a city. However, from the past performance of many electric bulbs of a certain make, or from the meteorological records over many years in the past, we may be able to assign probabilities to various outcomes in these two examples. Such probabilities are empirical in nature.

## 1.2 Set Theory

The collection of all possible outcomes of a random experiment and various sub-collections of these outcomes are the entities to which we want to assign probabilities.

## 2 THEORY AND METHODS OF STATISTICS

These sub-collections are *sets* in the *space* of all possible outcomes. Our aim is to develop a logical system which would enable us to calculate the probabilities of sets of complicated nature from the probabilities of sets whose probabilities are more clearly understood. For this, we need to have an understanding of some basic facts about the theory of sets.

We denote by  $S$  the space of all possible outcomes of a random experiment, consisting of elements (particular outcomes)  $s \in S$ . In  $S$ , we would be interested in various sets  $A \subset S$  and their combinations of different types.

### Definition 1.2.1.

- (i) The entire space is  $S$  consisting of elements  $s \in S$ , and  $\emptyset$  is the empty set which contains no element of  $S$ .
- (ii)  $A^c = \{s \in S: s \notin A\}$  is the complement of  $A \subset S$ .
- (iii)  $A_1 \cap A_2 = A_1 A_2 = \{s \in S: s \in A_1 \text{ and } s \in A_2\}$  is the intersection of  $A_1$  and  $A_2$ .
- (iv)  $A_1 \cup A_2 = \{s \in S: s \in \text{at least one of } A_1, A_2\}$  is the union of  $A_1$  and  $A_2$ .
- (v)  $A_1 \subset A_2$  if  $s \in A_1$  implies  $s \in A_2$ ;  $A_1 = A_2$  if  $A_1 \subset A_2$  and  $A_2 \subset A_1$ .

The following are immediate from the above definitions:

- (i)  $\emptyset = S^c$ .
- (ii)  $A_1 \cup A_2 = (A_1 A_2^c) \cup (A_1^c A_2) \cup (A_1 A_2)$ .
- (iii)  $(A^c)^c = A$ .
- (iv)  $(A_1 \cup A_2)^c = A_1^c \cap A_2^c$  and more generally  $(\bigcup_{i=1}^n A_i)^c = \bigcap_{i=1}^n A_i^c$ .
- (v)  $(A_1 \cap A_2)^c = A_1^c \cup A_2^c$  and more generally  $(\bigcap_{i=1}^n A_i)^c = \bigcup_{i=1}^n A_i^c$ .
- (vi)  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$  and  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ .
- (vii) If  $A_1 \subset A_2$ , then  $A_1 \cup A_2 = A_2$  and  $A_1 \cap A_2 = A_1$ .

**Note.** (iv) and (v) are called DeMorgan's rules.

**Definition 1.2.2.** Sets  $A_1, A_2, \dots$  are disjoint (or mutually exclusive) if  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

### Definition 1.2.3.

- (i) If the sets  $A_1, A_2, \dots$  are such that  $\bigcup_{n=1}^{\infty} A_n = S$ , then the collection  $\{A_1, A_2, \dots\}$  forms a covering of  $S$ .
- (ii) If the sets  $A_1, A_2, \dots$  are disjoint and  $\bigcup_{n=1}^{\infty} A_n = S$ , then the collection  $\{A_1, A_2, \dots\}$  forms a partition of  $S$ .

**Note.** To avoid triviality, the sets in a covering or a partition of  $S$  are nonempty.

From any covering  $\{A_1, A_2, \dots\}$  of  $S$  we can construct a partition  $\{B_1, B_2, \dots\}$  of  $S$  by letting  $B_1 = A_1, B_2 = A_1^c A_2, \dots, B_n = A_1^c A_2^c \cdots A_{n-1}^c A_n, \dots$

Clearly,  $B_1, B_2, \dots$  are disjoint and  $\bigcup_{n=1}^{\infty} B_n = \bigcup_{n=1}^{\infty} A_n = S$ . In particular, if  $A_1 \subset A_2 \subset \cdots \subset A_n \subset \cdots$  and  $\bigcup_{n=1}^{\infty} A_n = S$ , then  $\{B_1, B_2, \dots\}$  forms a partition of  $S$  if  $B_1 = A_1, B_2 = A_1^c A_2, \dots, B_n = A_{n-1}^c A_n, \dots$

## 1.3 Axiomatic Definition of Probability

A *probability space* is described by a triple  $(S, \mathcal{A}, P)$  where  $S$  is an arbitrary space consisting of elements  $s \in S$ ,  $\mathcal{A}$  is a collection of sets  $A \subset S$ , called *events*, with the properties:

$$(Ei) S \in \mathcal{A}. (Eii) A \in \mathcal{A} \Rightarrow A^c \in \mathcal{A}. (Eiii) A_1, A_2, \dots \in \mathcal{A} \Rightarrow \bigcup_{n=1}^{\infty} A_n \in \mathcal{A},$$

and  $P: \mathcal{A} \rightarrow [0, 1]$  is a function on  $\mathcal{A}$  with the properties:

$$(Pi) P[S] = 1. (Pii) \text{ for disjoint sets } A_1, A_2, \dots \in \mathcal{A}, P\left[\bigcup_{n=1}^{\infty} A_n\right] = \sum_{n=1}^{\infty} P[A_n].$$

A collection of sets of  $S$  with properties E(i, ii, iii) is called a  *$\sigma$ -field of subsets of  $S$*  and a function on  $\mathcal{A}$  with the property P(ii) is called a *a countably additive set function*. The property P(i) makes such a set function a *Probability*.

**Proposition 1.3.1** (Continuity Property of Probability). *From the axioms P(i, ii), it follows that*

- (Piii) if  $A_1 \supset A_2 \supset A_3 \supset \dots \supset A_n \supset \dots$  and  $\lim_{n \rightarrow \infty} A_n = \bigcap_{n=1}^{\infty} A_n = \emptyset$ , then  $\lim_{n \rightarrow \infty} P[A_n] = 0$ , or equivalently,
- (Piv) if  $A_1 \subset A_2 \subset A_3 \subset \dots \subset A_n \subset \dots$  and  $\lim_{n \rightarrow \infty} A_n = \bigcup_{n=1}^{\infty} A_n = S$ , then  $\lim_{n \rightarrow \infty} P[A_n] = 1$ .

Conversely, either P(i, iii) or P(i, iv) implies P(i, ii).

Proof of the equivalence of P(i, ii), P(i, iii), and P(i, iv) is left as an exercise.

## 1.4 Some Simple Propositions

### Proposition 1.4.1.

- (i)  $P[A^c] = 1 - P[A]$ ,  $P[\emptyset] = 0$ .
- (ii) If  $A \subset B$ , then  $P[A] \leq P[B]$ .
- (iii)  $P[A \cup B] = P[A] + P[B] - P[AB]$ .
- (iv)  $P\left[\bigcup_{i=1}^n A_i\right] \leq \sum_{i=1}^n P[A_i]$ .

*Proof.*

- (i)  $A$  and  $A^c$  are disjoint and  $A \cup A^c = S$ . Hence  $1 = P[S] = P[A \cup A^c] = P[A] + P[A^c]$ , so  $P[A^c] = 1 - P[A]$ . In particular,  $\emptyset = S^c$  and therefore,  $P[\emptyset] = 1 - P[S] = 1 - 1 = 0$ .
- (ii) If  $A \subset B$ , then  $B = A \cup (A^c B)$ , where  $A$  and  $A^c B$  are disjoint. Hence  $P[B] = P[A] + P[A^c B] \geq P[A]$ , because  $P[A^c B] \geq 0$ .
- (iii)  $A = (AB) \cup (AB^c)$ ,  $B = (AB) \cup (A^c B)$ ,  $A \cup B = (AB) \cup (A^c B) \cup (AB)$ , and the sets  $AB^c$ ,  $A^c B$ ,  $AB$  are disjoint. Thus

$$\begin{aligned} P[A \cup B] &= P[AB^c] + P[A^c B] + P[AB] \\ &= \{P[AB] + P[AB^c]\} + \{P[AB] + P[A^c B]\} - P[AB] \\ &= P[A] + P[B] - P[AB]. \end{aligned}$$

## 4 THEORY AND METHODS OF STATISTICS

- (iv) Let  $B_1 = A_1, B_2 = A_1^c A_2, \dots, B_n = A_1^c \cdots A_{n-1}^c A_n$ . Then
- (a)  $\bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i$ ,
  - (b)  $B_1, B_2, \dots, B_n$  are disjoint, and
  - (c)  $B_i \subset A_i, i = 1, \dots, n$ .
- Hence  $P\left[\bigcup_{i=1}^n A_i\right] = P\left[\bigcup_{i=1}^n B_i\right] = \sum_{i=1}^n P[B_i] \leq \sum_{i=1}^n P[A_i]$  by (ii).

□

### Proposition 1.4.2.

$$\begin{aligned} P[A_1 \cup \dots \cup A_n] &= \sum_{i=1}^n P[A_i] - \sum_{1 \leq i_1 < i_2 \leq n} P[A_{i_1} A_{i_2}] \\ &\quad + \dots + (-1)^{r+1} \sum_{1 \leq i_1 < i_2 < \dots < i_r \leq n} P[A_{i_1} A_{i_2} \cdots A_{i_r}] \\ &\quad + \dots + (-1)^{n+1} P[A_1 A_2 \cdots A_n] \end{aligned}$$

The proof, which follows by induction, starting from [Proposition 1.4.1](#)(iii), is left as an exercise.

**Proposition 1.4.3** (Bonferroni's Inequality).  $P[A_1 A_2 \cdots A_n] \geq P[A_1] + P[A_2] + \dots + P[A_n] - 1$ .

For  $n = 2$ ,  $P[A_1 A_2] \geq P[A_1] + P[A_2] - 1$  is proved as follows:  $1 \geq P[A_1 \cup A_2] = P[A_1] + P[A_2] - P[A_1 A_2]$ . Hence the result. The proof of the general case follows by induction, and is left as an exercise.

This inequality may be useful in some situations to get some idea about  $P[A_1 A_2 \cdots A_n]$  where actual evaluation may be difficult.

## 1.5 Equally Likely Outcomes in Finite Sample Space

Suppose that  $S$  consists of  $N$  elements,  $S = \{s_1, \dots, s_N\}$ , which are equally likely. Then the total probability of 1 is shared equally by these  $N$  elements, so that  $P[\{s_i\}] = \frac{1}{N}$  for each  $i$ . Consequently, for any event  $A = \{s_{i_1}, \dots, s_{i_K}\}$ ,  $P[A] = \frac{K}{N} = \frac{\#(A)}{\#(S)}$ , where  $\#(E)$  ≈ number of elements in  $E$ .

Evaluation of probabilities of events in such a setting, essentially reduces to counting problems which are often of a combinatorial nature.

**Example 1.5.1.** Two balanced dice with 1, 2, 3, 4, 5, 6 on the six faces are rolled. Find the probabilities of the events: (a) the numbers on the two dice are equal, (b) the sum of the numbers on the two dice is 10, (c) the number on at least one die is 6, and (d) none of the numbers is a 6.

*Solution.*

(a)  $\#(S) = 6 \cdot 6 = 36$  and if we let  $A_1$  be the event in question, then

$A_1 = \{(1, 1), (2, 2), \dots, (6, 6)\}$  and  $\#(A_1) = 6$ . Thus  $P[A_1] = \frac{\#(A_1)}{\#(S)} = \frac{6}{36} = \frac{1}{6}$ .

(b) Here the event  $A_2 = \{(4, 6), (5, 5), (6, 4)\}$ , so  $\#(A_2) = 3$ . Thus  $P[A_2] = \frac{\#(A_2)}{\#(S)} = \frac{3}{36} = \frac{1}{12}$ .

- (d) We answer (d) first. None of the numbers is a 6, means no 6 on the first, with probability  $\left(1 - \frac{1}{6}\right)$  and no 6 on the second, with probability  $\left(1 - \frac{1}{6}\right)$ . Hence the probability of no 6 on any of the two dice is the probability of the event  $A_4$ ,  $P[A_4] = (1 - \frac{1}{6})(1 - \frac{1}{6}) = \frac{5}{6} \cdot \frac{5}{6} = \frac{25}{36}$ , because the roll of the first die has nothing to do with the roll of the second die, that is, they are “independent,” which is a concept to be introduced in the sequel. In a direct counting approach,  $A_4 = \{(x, y) : x = 1, 2, 3, 4, 5 \text{ and } y = 1, 2, 3, 4, 5\}$ . Thus  $\#(A_4) = 5 \cdot 5 = 25$ , so  $P[A_4] = \frac{\#(A_4)}{\#(S)} = \frac{25}{36}$ , which agrees with our calculation based on the concept of independence.
- (c) “At least one 6” is the complement of “no 6”. Hence  $P[A_3] = P[A_4^C] = 1 - P[A_4] = 1 - \frac{25}{36} = \frac{11}{36}$ .

**Example 1.5.2.** What is the probability that a bridge hand dealt from a well-shuffled standard deck of cards will contain 2 Aces and 2 Kings?

*Solution.* Here  $\#(S) \doteq$  number of ways to choose 13 cards from a deck of 52 cards  $= \binom{52}{13}$ , and if  $A$  denotes the event in question, then  $\#(A) \doteq$  number of ways to choose 2 out of 4 Aces, 2 out of 4 Kings, and 9 out of the remaining 44 cards  $= \binom{4}{2} \binom{4}{2} \binom{44}{9}$ . Thus

$$P[A] = \frac{\binom{4}{2} \binom{4}{2} \binom{44}{9}}{\binom{52}{13}} \approx 0.04.$$

## 1.6 Conditional Probability and Independence

**Definition 1.6.1.** For events  $A$  and  $B$  with  $P[A] > 0$ , the conditional probability of  $B$  given  $A$  is defined to be

$$P[B|A] = \frac{P[AB]}{P[A]}. \quad (1)$$

This is the probability that we assign to the event  $B$  if we know that the event  $A$  has occurred. Now if the knowledge of occurrence of  $A$  does not change the probability of  $B$  (ie, if  $P[B|A] = P[B]$ ), then we say that  $B$  is *independent of A*. But this holds if  $P[AB] = P[A]P[B]$ . Now if  $P[B] > 0$ , then  $P[AB] = P[A]P[B]$  also implies

$$P[A|B] = \frac{P[AB]}{P[B]} = P[A],$$

so that  $A$  is *independent of B*. This leads to the following definition.

**Definition 1.6.2.** Events  $A, B$  are independent if

$$P[AB] = P[A]P[B]. \quad (2a)$$

More generally, events  $A_1, \dots, A_n$  are independent if for every subset  $\{i_1, \dots, i_r\}$  of  $\{1, 2, \dots, n\}$ ,  $2 \leq r \leq n$ ,

$$P[A_{i_1} A_{i_2} \cdots A_{i_r}] = P[A_{i_1}]P[A_{i_2}] \cdots P[A_{i_r}]. \quad (2b)$$

## 6 THEORY AND METHODS OF STATISTICS

From the definition of  $P[B|A]$  in Eq. (1), we have

$$P[AB] = P[A]P[B|A]. \quad (3a)$$

More generally, we have

$$P[A_1A_2 \cdots A_n] = P[A_1]P[A_2|A_1]P[A_3|A_1A_2] \cdots P[A_n|A_1 \cdots A_{n-1}]. \quad (3b)$$

This is proved by applying the above formula repeatedly to the right side to see that

$$\begin{aligned} & P[A_1]P[A_2|A_1]P[A_3|A_1A_2] \cdots P[A_n|A_1 \cdots A_{n-1}] \\ &= P[A_1A_2]P[A_3|A_1A_2] \cdots P[A_n|A_1 \cdots A_{n-1}] \\ &= P[A_1A_2A_3] \cdots P[A_n|A_1 \cdots A_{n-1}] \\ &= \cdots = P[A_1A_2 \cdots A_n]. \end{aligned}$$

Formulas (2a) and (2b) are special cases of Eqs. (3a) and (3b) when the events are independent.

**Proposition 1.6.1.** [Bayes Formula] If the collection of events  $\{A_1, A_2, \dots\}$  forms a partition of  $S$ , that is, if  $A_1, A_2, \dots \in \mathcal{A}$  are disjoint and  $\bigcup_{j=1}^{\infty} A_j = S$ , then

$$P[A_i|B] = \frac{P[A_iB]}{P[B]} = \frac{P[A_i]P[B|A_i]}{\sum_{j=1}^{\infty} P[A_j]P[B|A_j]}.$$

*Proof.* We only have to observe that by virtue of  $A_1, A_2, \dots$  being a partition of  $S$ ,  $B = \bigcup_{j=1}^{\infty} A_j B$ , and the events  $A_1B, A_2B, \dots$  are disjoint. Hence  $P[B] = \sum_{j=1}^{\infty} P[A_jB] = \sum_{j=1}^{\infty} P[A_j]P[B|A_j]$  by P(ii) and Eq. (3a).  $\square$

For any event  $B$  with  $P[B] > 0$ ,  $P[A|B]$  as a function of  $A \in \mathcal{A}$  is a probability, that is,  $0 \leq P[A|B] \leq 1$  for all  $A \in \mathcal{A}$  and  $P[\cdot|B]$  satisfies P(i, ii). The verification of these facts and the proof of the following proposition are left as exercises.

**Proposition 1.6.2.** If  $A$  and  $B$  are independent, then the pairs  $A$  and  $B^c$ ,  $A^c$  and  $B$ ,  $A^c$  and  $B^c$  are also independent.

**Example 1.6.1.** There are 20 balls in a bag, of which 8 are white and 12 are black. We draw three balls at random from the bag. Find the probability that all three of these balls are white if (a) each ball drawn is replaced before the next draw, (b) the balls are drawn without replacement.

*Solution.*

- (a) Here the outcomes of the three draws are independent events, because the composition of the bag is unchanged after a draw due to replacement. Let  $A_i$  denote the event of a white ball in  $i$ th draw,  $i = 1, 2, 3$ . Then  $P[A_i] = \frac{8}{20}$  for each  $i$  and  $A_1, A_2, A_3$  are independent. Hence  $P(\text{All three balls are white}) = P[A_1A_2A_3] = P[A_1]P[A_2]P[A_3] = \left(\frac{8}{20}\right)^3 = \left(\frac{2}{5}\right)^3 = \frac{8}{125}$ .
- (b) When the draws are without replacement, the events  $A_1, A_2, A_3$  are not independent. Here we can approach the problem in two ways:

(i) By direct counting, we have  $\#(A_1 A_2 A_3) = \binom{8}{3} = 56$  and  $\#(S) = \binom{20}{3} = 1140$ . Hence  $P[A_1 A_2 A_3] = \frac{56}{1140} = \frac{14}{285}$ .

(ii) Arguing with conditional probabilities, by looking at the composition of the bag conditional upon the outcomes of previous draws we have

$$P[A_1] = \frac{8}{20}, P[A_2|A_1] = \frac{7}{19}, P[A_3|A_1 A_2] = \frac{6}{18}.$$

$$\text{Hence } P[A_1 A_2 A_3] = P[A_1] P[A_2|A_1] P[A_3|A_1 A_2] = \frac{8}{20} \cdot \frac{7}{19} \cdot \frac{6}{18} = \frac{14}{285}.$$

**Example 1.6.2.** There are three boxes of which Box 1 contains tickets numbered  $1, \dots, 5$ ; Box 2 contains tickets numbered  $1, \dots, 10$ ; and Box 3 contains tickets numbered  $1, \dots, 20$ . One of the three boxes is chosen at random and then a ticket is drawn at random from the chosen box. Find the probability that the ticket is from Box 2 if (a) the number on the ticket is 4, (b) the number on the ticket is 7.

*Solution.* Let  $A_i$  be the event that Box  $i$  is chosen,  $B$  the event that the number on the ticket is 4 and  $C$  the event that the number on the ticket is 7. Then  $P[A_i] = \frac{1}{3}, i = 1, 2, 3$ ;  $P[B|A_1] = \frac{1}{5}, P[B|A_2] = \frac{1}{10}, P[B|A_3] = \frac{1}{20}$ ;  $P[C|A_1] = 0, P[C|A_2] = \frac{1}{10}, P[C|A_3] = \frac{1}{20}$ .

(a) Here we need  $P[A_2|B]$ . By Bayes' formula,

$$P[A_2|B] = \frac{P[A_2]P[B|A_2]}{\sum_{j=1}^3 P[A_j]P[B|A_j]} = \frac{\frac{1}{3} \cdot \frac{1}{10}}{\frac{1}{3} \cdot \frac{1}{5} + \frac{1}{3} \cdot \frac{1}{10} + \frac{1}{3} \cdot \frac{1}{20}} = \frac{2}{7}.$$

(b) Here we need  $P[A_2|C]$ . By Bayes' formula,

$$P[A_2|C] = \frac{P[A_2]P[C|A_2]}{\sum_{j=1}^3 P[A_j]P[C|A_j]} = \frac{\frac{1}{3} \cdot \frac{1}{10}}{\frac{1}{3} \cdot 0 + \frac{1}{3} \cdot \frac{1}{10} + \frac{1}{3} \cdot \frac{1}{20}} = \frac{2}{3}.$$

## 1.7 Random Variables and Their Distributions

**Definition 1.7.1.** A real-valued function  $X: S \rightarrow \mathbb{R}$  on a probability space  $(S, \mathcal{A}, P)$  is a random variable (*rv*) if for all  $a \in \mathbb{R}$ ,

$$\{s \in S: X(s) \leq a\} \in \mathcal{A}. \quad (4a)$$

More generally,  $\mathbf{X} = (X_1, \dots, X_k): S \rightarrow \mathbb{R}^k$  on a probability space  $(S, \mathcal{A}, P)$  is a  $k$ -dimensional ( $k$ -dim) random vector (*rv*) if for all  $(a_1, \dots, a_k) \in \mathbb{R}^k$ ,

$$\{s \in S: X_1(s) \leq a_1, \dots, X_k(s) \leq a_k\} \in \mathcal{A}. \quad (4b)$$

**Definition 1.7.2.** The cumulative distribution function (cdf) of a 1-dim rv  $X$  is defined as

$$F_X(a) = P[X \leq a] = P[\{s: X(s) \leq a\}] \quad \text{for all } a \in \mathbb{R}. \quad (5a)$$

**Note.** We shall use  $P$  as a set function  $P: \mathcal{A} \rightarrow [0, 1]$  as in P(i, ii) as well as in a more informal way as in  $P[X \leq a]$  to indicate the probability of  $X \leq a$ .

## 8 THEORY AND METHODS OF STATISTICS

*Remark 1.7.1.* We are using the abbreviation rv to indicate both a random variable and a random vector. The meaning of one or the other should be clear in the context.

**Proposition 1.7.1** (Properties of a cdf  $F_X(x)$ ).

- (i)  $F_X(a) \leq F_X(b)$  for all  $a \leq b$ .
- (ii)  $\lim_{h \downarrow 0} F_X(x+h) = F_X(x)$  for all  $x$ .
- (iii)  $\lim_{x \rightarrow -\infty} F_X(x) = 0$ .
- (iv)  $\lim_{x \rightarrow \infty} F_X(x) = 1$ .

Property (i) says that  $F_X$  is monotone nondecreasing (because  $(-\infty, a] \subset (-\infty, b]$  for all  $a \leq b$ ) and Property (ii) says that  $F_X$  is right continuous. Properties (ii), (iii), and (iv) follow from the continuity property of probability, the proofs of which are left as exercises.

For any rv  $X$ ,

$$\begin{aligned} P[X = a] &= P[X \leq a] - P[X < a] = F_X(a) - \lim_{h \uparrow 0} F_X(a+h) \\ &= \text{magnitude of jump of } F_X \text{ at } a \text{ (if any).} \end{aligned}$$

We shall consider only two types of rv's, discrete and continuous. A rv  $X$  is *discrete* if its cdf  $F_X$  is a step-function (ie, increases only by jumps and stays constant elsewhere) and is *continuous* if its cdf  $F_X$  is differentiable everywhere. Since a cdf  $F_X$  can have only a finite or a countably infinite number of jumps, a discrete rv can have only a finite or countably infinite number of possible values  $\{x_i\}$ .

**Definition 1.7.3.** The probability mass function (pmf) of a discrete rv  $X$  is defined as

$$f_X(x_i) = P[X = x_i], \quad i = 1, 2, \dots \quad \text{and} \quad f_X(x) = 0 \quad \text{for all } x \notin \{x_i\}. \quad (6a)$$

For any set  $B \subset \mathbb{R}$ ,  $P[X \in B] = \sum_{x \in B} f_X(x)$  and  $\sum_{i=1}^{\infty} f_X(x_i) = P[X \in \mathbb{R}] = 1$ .

**Definition 1.7.4.** The probability density function (pdf) of a continuous rv  $X$  is defined as

$$f_X(x) = F'_X(x) \geq 0 \quad \text{for all } x \in \mathbb{R}. \quad (7a)$$

Clearly,  $P[X \leq a] = F_X(a) = \int_{-\infty}^a f_X(x) dx$ , and  $\int_{-\infty}^{\infty} f_X(x) dx = \lim_{x \rightarrow \infty} F_X(x) = 1$ .

**Note.** Since the cdf of a continuous rv  $X$  is differentiable everywhere, it has no jumps. Hence  $P[X = a] = 0$  for all  $a$ .

On the real line  $\mathbb{R}$ , the smallest  $\sigma$ -field of subsets which includes all intervals, is called the  $\sigma$ -field of Borel sets and is denoted by  $\mathbb{B}$ . For each  $B \in \mathbb{B}$ ,

$$P[X \in B] = \int_B f_X(x) dx,$$

when this integral is defined in an appropriate manner.

The cdf of a  $k$ -dim rv  $\mathbf{X} = (X_1, \dots, X_k)$  is defined as

$$F_{\mathbf{X}}(a_1, \dots, a_k) = P[X_1 \leq a_1, \dots, X_k \leq a_k] \quad \text{for all } (a_1, \dots, a_k) \in \mathbb{R}^k. \quad (5b)$$

The marginal cdf of  $X_i$  is obtained as  $F_{X_i}(a) = F_{\mathbf{X}}(\infty, \dots, \infty, a, \infty, \dots, \infty)$ , that is, by taking  $a_j = \infty$ , and  $a_j = a$  for all  $j \neq i$  in  $F_{\mathbf{X}}(a_1, \dots, a_k)$ . Again we consider the discrete and continuous cases, and define

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \begin{cases} P[X_1 = x_1, \dots, X_k = x_k] & \text{if } (x_1, \dots, x_k) \in \text{set of possible values of } \mathbf{X} \\ 0 & \text{otherwise} \end{cases} \quad (6b)$$

as the pmf of  $\mathbf{X}$  in the discrete case, and

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \frac{\partial^k}{\partial x_1 \cdots \partial x_k} F_{\mathbf{X}}(x_1, \dots, x_k) \quad (7b)$$

as the pdf of  $\mathbf{X}$  in the continuous case.

The various marginal pmf's and pdf's are obtained by summing or integrating over the other variables in Eq. (6b) or (7b).

For random variables defined by (4a, b), the cdf defined by (5a, b), the pmf/pdf defined by (6a, b) or (7a, b) contain all the information.

**Example 1.7.1.** Players 1 and 2 each rolls a balanced six-faced die. Let  $(s_1, s_2)$  be the outcome of the rolls,  $s_i$  the number on the die rolled by Player  $i$ . On the sample space of  $\{(s_1, s_2): s_i = 1, 2, 3, 4, 5, 6, i = 1, 2\}$ , we define  $X(s_1, s_2) = s_1 - s_2$  as the amount which Player 1 wins from Player 2, a negative win being interpreted as a loss. Find the pmf of  $X$ .

*Solution.* The possible values of the rv  $X$  are  $0, \pm 1, \pm 2, \pm 3, \pm 4$ , and  $\pm 5$ . For  $x \geq 0$ ,  $\#\{(s_1, s_2): s_1 - s_2 = x\} = \#\{(1+x, 1), \dots, (6, 6-x)\} = 6 - x$ , each carrying a probability of  $\frac{1}{36}$ . Hence for  $x \geq 0$ ,

$$f_X(x) = P[X = x] = \frac{6-x}{36}, \quad x = 0, 1, \dots, 5.$$

Again for  $x < 0$ ,  $\#\{(s_1, s_2): s_1 - s_2 = x\} = \#\{(s_1, s_2): s_2 - s_1 = |x|\} = \#\{(1, 1+|x|), \dots, (6-|x|, 6)\} = 6 - |x|$ , each carrying a probability of  $\frac{1}{36}$ . Hence for  $x < 0$ ,

$$f_X(x) = P[X = x] = \frac{6-|x|}{36}, \quad x = -5, \dots, -1.$$

Thus,  $f_X(x) = \frac{6-|x|}{36}$ ,  $x = 0, \pm 1, \pm 2, \pm 3, \pm 4, \pm 5$ .

**Example 1.7.2.** A point  $s$  is selected at random from the unit circle with center at the origin  $(0, 0)$  and radius 1 and

$$X(s) = \text{distance of } s \text{ from the center.}$$

Find the cdf and the pdf of  $X$  and the probability that  $\frac{1}{2} \leq X \leq \frac{3}{4}$ .

*Solution.* The cdf of  $X$  is

$$\begin{aligned} F_X(x) &= P[\{s: X(s) \leq x\}] \\ &= P[s \in \text{a circle of radius } x \text{ centered at origin}] \\ &= \frac{\text{Area of a circle of radius } x}{\text{Area of unit circle}} \\ &= \frac{\pi x^2}{\pi 1^2} = x^2, \quad 0 \leq x \leq 1. \end{aligned}$$

This is because, for a randomly selected point  $s$  from a set  $S$  in the plane,

$$P[s \in A] = \frac{\text{Area of } A}{\text{Area of } S}, \quad \text{for all } A \subset S$$

assuming that the areas of  $S$  and of  $A$  are well-defined.

Now the pdf of  $X$  is

$$f_X(x) = F'_X(x) = \frac{d}{dx}(x^2) = 2x, \quad 0 \leq x \leq 1.$$

Finally,

$$P\left[\frac{1}{2} \leq X \leq \frac{3}{4}\right] = F_X\left(\frac{3}{4}\right) - F_X\left(\frac{1}{2}\right) = \left(\frac{3}{4}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{5}{16}.$$

Alternatively,

$$P\left[\frac{1}{2} \leq X \leq \frac{3}{4}\right] = \int_{\frac{1}{2}}^{\frac{3}{4}} f_X(x) dx = \int_{\frac{1}{2}}^{\frac{3}{4}} 2x dx = x^2 \Big|_{\frac{1}{2}}^{\frac{3}{4}} = \frac{5}{16}.$$

## 1.8 Expected Value, Variance, Covariance, and Correlation Coefficient

**Definition 1.8.1.** The expected value or the mean of an rv  $X$  is defined as

$$\mu_X = E[X] = \begin{cases} \sum_{i=1}^{\infty} x_i f_X(x_i) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} x f_X(x) dx & \text{if } X \text{ is continuous,} \end{cases}$$

provided that the sum or the integral exists and is finite. Otherwise, we say that  $E[X]$  does not exist.

More generally, we define

$$E[g(X)] = \begin{cases} \sum_{i=1}^{\infty} g(x_i) f_X(x_i) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} g(x) f_X(x) dx & \text{if } X \text{ is continuous.} \end{cases}$$

**Note.**  $E[X]$  exists if and only if  $E[|X|] < \infty$ .

For a  $k$ -dim random vector  $\mathbf{X} = (X_1, \dots, X_k)$ , we define

$$\begin{aligned} E[g(X_1, \dots, X_k)] \\ = \begin{cases} \sum_{\text{all } x_1, \dots, x_k} g(x_1, \dots, x_k) f_X(x_1, \dots, x_k) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1, \dots, x_k) f_X(x_1, \dots, x_k) dx_1 \cdots dx_k & \text{if } X \text{ is continuous.} \end{cases} \end{aligned}$$

**Note.** The term *Expectation* is often used for *Expected Value*.

**Proposition 1.8.1** (Properties of Expectation).

- (i) If  $P[X = c] = 1$ , then  $E[X] = c$ .
- (ii) If  $P[X \geq 0] = 1$  with probability 1, then  $E[X] \geq 0$ .
- (iii) If  $I_A(x) = 1$  for  $x \in A$  and  $I_A(x) = 0$  for  $x \notin A$ , then  $E[I_A(x)] = P[X \in A]$ .
- (iv)  $E[aX + b] = aE[X] + b$ .
- (v)  $E\left[\sum_{i=1}^k a_i X_i\right] = \sum_{i=1}^k a_i E[X_i]$ .
- (vi) If  $P[X_1 \geq X_2] = 1$ , then  $E[X_1] \geq E[X_2]$ .

The proof is omitted.

**Definition 1.8.2.** The variance of an rv  $X$  is defined as

$$\sigma_X^2 = \text{Var}[X] = E[(X - E[X])^2] = E[X^2] - \{E[X]\}^2.$$

The last equality follows by expanding  $(X - E[X])^2$  and then using the properties of expectation in [Proposition 1.8.1](#), remembering that  $E[X] = \mu_X$  is a constant. The *standard deviation* of  $X$  is  $\sqrt{\text{Var}[X]} = \sigma_X$ .

**Definition 1.8.3.** The covariance between rv's  $X$  and  $Y$  is defined as

$$\sigma_{XY} = \text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y].$$

Again the last equality follows by expanding  $(X - E[X])(Y - E[Y])$  and using the properties of expectation. Obviously,

$$\text{Cov}[X, X] = \text{Var}[X].$$

**Note.** If  $\mathbf{X}^\top = (X_1, \dots, X_k)$  is a  $k$ -dim rv, then  $\boldsymbol{\Sigma} = ((\text{Cov}[X_i, X_j])_{i,j})$  is called the *covariance matrix of  $X$* .

**Proposition 1.8.2** (Properties of Variance and Covariance).

- (i)  $\text{Var}[aX + b] = a^2 \sigma_X^2$ .
- (ii)  $\text{Var}\left[\sum_{i=1}^k a_i X_i\right] = \sum_{i=1}^k a_i^2 \sigma_{X_i}^2 + 2 \sum_{1 \leq i < j \leq k} a_i a_j \sigma_{X_i X_j}$ .
- (iii)  $\text{Cov}[aX + b, cY + d] = ac \sigma_{X,Y}$ .
- (iv)  $\text{Cov}\left[\sum_{i=1}^m a_i X_i + \sum_{j=1}^n b_j Y_j\right] = \sum_{i=1}^m \sum_{j=1}^n a_i b_j \sigma_{X_i Y_j} = \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{b}$ ,  
where  $\mathbf{a}^\top = (a_1, \dots, a_m)$ ,  $\mathbf{b}^\top = (b_1, \dots, b_n)$ , and  $\boldsymbol{\Sigma} = ((\sigma_{(X_i, Y_j)}))$ .
- (v) If  $X$  is a  $k$ -dim rv, then  $\text{Cov}[X, X] = E[(X - \mu_X)(X - \mu_X)^\top]$ .
- (vi) If  $X$  is a  $k$ -dim rv and  $A$  is a  $k \times k$  matrix, then  $\text{Cov}[AX, AX] = AC\text{Cov}[X, X]A^\top$ .

The proof is omitted.

**Question.** Why is the mean  $\mu_X$  called the “expected value” of  $X$ ?

The answer is statistical in nature. Suppose we want to *predict* the value of the rv  $X$  before it is actually observed, by making a guess  $a$  and paying a penalty  $(X - a)^2$ , which is the square of the discrepancy between our guess  $a$  and the actual observation  $X$ . Then the average penalty is

$$\begin{aligned} E[(X - a)^2] &= E[(X - \mu_X) - (a - \mu_X)]^2 \\ &= E[(X - \mu_X)^2] + (a - \mu_X)^2 = \sigma_X^2 + (a - \mu_X)^2, \end{aligned}$$

(using the fact  $E[X - \mu_X] = 0$ ), which is minimized by taking  $a = \mu_X$ . Thus the mean  $\mu_X$  is the best guess for an unobserved  $X$  under the rule of squared-error penalty and that is why  $\mu_X = E[X]$  is called the expected value of  $X$ . In a variation of the above game of prediction, what would be the best guess if the penalty is the absolute error  $|X - a|$ ? Here we want to minimize  $E[|X - a|]$  by our choice of  $a$ . For simplicity, suppose  $F_X$  is continuous and strictly increasing and let  $m = m_X = F_X^{-1}(\frac{1}{2})$  be the solution of the equation  $F_X(m) = \frac{1}{2}$ . Then under the rule of absolute error penalty,  $m$  is the best guess for an unobserved  $X$ . The quantity  $m_X$  is called the *median* of  $X$ , which has some advantages over the mean  $\mu_X$ , as will be seen later. For one thing,  $m_X$  always exists even when  $\mu_X$  may not.

We now prove an important inequality.

**Proposition 1.8.3** (Cauchy-Schwarz Inequality). *Let  $X, Y$  be rv's for which  $E[X^2] < \infty$  and  $E[Y^2] < \infty$ . Then*

$$\{E[XY]\}^2 \leq E[X^2]E[Y^2].$$

*Proof.* The function

$$h(t) = E[(tX - Y)^2] = t^2E[X^2] - 2tE[XY] + E[Y^2] \geq 0 \quad \text{for all } t,$$

so the quadratic equation  $h(t) = 0$  can have *at most one real root*. But the roots of this equation are

$$t = \frac{2E[XY] \pm \sqrt{4\{E[XY]\}^2 - 4E[X^2]E[Y^2]}}{2E[X^2]}.$$

Hence the expression under the radical must be  $\leq 0$ . □

Replacing  $X$  and  $Y$  by  $X - \mu_X$  and  $Y - \mu_Y$ , respectively, in the above inequality, we have

$$\text{Cov}^2[X, Y] \leq \text{Var}[X] \text{Var}[Y], \quad \text{ie, } \sigma_{XY}^2 \leq \sigma_X^2 \sigma_Y^2.$$

**Definition 1.8.4.** The correlation coefficient between  $X$  and  $Y$  is defined as

$$\rho_{XY} = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]\text{Var}[Y]}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}.$$

From *Cauchy-Schwarz inequality*, it follows that  $\rho_{XY}^2 \leq 1$ , that is,  $-1 \leq \rho_{XY} \leq 1$ . Moreover,  $\rho_{XY} = +1$  or  $\rho_{XY} = -1$ , that is,  $\rho_{XY}^2 = 1$  iff the equation

$$h(t) = E[t(X - \mu_X) - (Y - \mu_Y)]^2 = 0$$

holds for exactly one real value of  $t$ . But

$$\begin{aligned} E[t(X - \mu_X) - (Y - \mu_Y)]^2 &= 0 \\ \Leftrightarrow Y - \mu_Y &= t(X - \mu_X) \quad \text{w.p. 1 for some } t \neq 0. \end{aligned}$$

( $t = 0 \Rightarrow Y - \mu_Y = 0$  w.p. 1  $\Rightarrow \sigma_Y^2 = 0$ , and we exclude this degenerate case from our consideration.) Thus  $\rho_{XY} = +1$  or  $-1$  iff  $Y$  is a linear function of  $X$  w.p. 1, the slope having the same sign as  $\rho_{XY}$ .

## 1.9 Moments and the Moment Generating Function

**Definition 1.9.1.** The  $r$ th moment of an rv  $X$  is defined as

$$\mu_{r,X} = E[X^r], \quad r = 1, 2, \dots$$

assuming existence, and  $\mu_{0,X} = 1$ . In particular, the mean  $\mu_X = \mu_{1,X}$ , and  $\sigma_X^2 = \mu_{2,X} - \mu_{1,X}^2$ .

**Definition 1.9.2.** The moment generating function (mgf) of  $X$  is defined as  $M_X(t) = E[e^{tX}]$ , provided that the expectation exists.

Differentiating under the integral,

$$\begin{aligned}\mu_X^{(r)}(0) &= \frac{d^r}{dt^r} M_X(t) \Big|_{t=0} = E[X^r e^{tX}] \Big|_{t=0} = E[X^r] = \mu_r, \\ M_X(t) &= \sum_{r=0}^{\infty} \frac{t^r}{r!} \mu_r.\end{aligned}$$

Hence the name moment generating function. For a linear function  $a + bX$ ,  $M_{a+bX}(t) = E[e^{t(a+bX)}] = e^{at} M_X(bt)$ .

## 1.10 Independent Random Variables and Conditioning When There Is Dependence

Random variables  $X_1, \dots, X_k$  are said to be mutually independent if

$$F_{X_1, \dots, X_k}(x_1, \dots, x_k) = \prod_{i=1}^k F_{X_i}(x_i) \quad \text{for all } x_1, \dots, x_k. \tag{8a}$$

Equivalently, for mutually independent rv's,

$$f_{X_1, \dots, X_k}(x_1, \dots, x_k) = \prod_{i=1}^k f_{X_i}(x_i) \quad \text{for all } x_1, \dots, x_k, \tag{8b}$$

holds for their joint pdf or joint pmf. The conditional pmf or pdf of  $X_{r+1}, \dots, X_k$  given  $(X_1, \dots, X_r) = (x_1, \dots, x_r)$  when  $f_{X_1, \dots, X_r}(x_1, \dots, x_r) > 0$  is

$$f_{(X_{r+1}, \dots, X_k)|(X_1, \dots, X_r)}(x_{r+1}, \dots, x_k | x_1, \dots, x_r) = \frac{f_{X_1, \dots, X_k}(x_1, \dots, x_k)}{f_{X_1, \dots, X_r}(x_1, \dots, x_r)}.$$

In particular, for  $(X, Y)$  having joint pmf/pdf  $f_{XY}(x, y)$ ,

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)} = \frac{f_{XY}(x, y)}{\int_{-\infty}^{\infty} f_{XY}(x, y) dy}, \quad \text{when } f_X(x) > 0.$$

On the other hand, from the marginal pdf  $f_X$  of  $X$  and the conditional pdf  $f_{Y|X}$  of  $Y$  given  $X$ , the joint pdf of  $(X, Y)$  is obtained as

$$f_{XY}(x, y) = f_X(x)f_{Y|X}(y|x). \quad (9)$$

If  $X_1, \dots, X_k$  are mutually independent rv's, then

$$\mathbb{E}[g_1(X_1) \cdots g_k(X_k)] = \prod_{i=1}^k \mathbb{E}[g_i(X_i)]$$

provided that  $\mathbb{E}[g_i(X_i)]$ ,  $i = 1, \dots, k$  exist. This follows immediately from the definition of independence (8a, b). It follows that if  $X, Y$  are independent, then  $\sigma_{XY} = \text{Cov}[X, Y] = 0$ , and  $\rho_{XY} = 0$ , and therefore,

$$\sigma_{aX+bY}^2 = a^2\sigma_X^2 + b^2\sigma_Y^2.$$

However,  $\rho_{XY} = 0$  does not imply that  $X, Y$  are independent.

Suppose that  $X_1, \dots, X_n$  are mutually independent and each  $X_i$  is distributed as  $X$ . We then say that  $X_1, \dots, X_n$  are independent and identically distributed (iid) as  $X$ , and if  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ , then

$$\mu_{\bar{X}} = \mu_X \text{ and } \sigma_{\bar{X}}^2 = \frac{\sigma_X^2}{n}.$$

If  $X_1, \dots, X_k$  are independent rv's and if the mgf  $M_{X_i}(t)$  exists for each  $X_i$ , then

$$M_{X_1+\dots+X_k}(t) = \mathbb{E}\left[e^{t \sum_{i=1}^k X_i}\right] = \prod_{i=1}^k \mathbb{E}[e^{tX_i}] = \prod_{i=1}^k M_{X_i}(t).$$

Going back to the general case of  $(X, Y)$  with joint pdf  $f_{XY}(x, y)$ , let

$$m(x) = \mathbb{E}[g(Y)|X = x] = \int g(y)f_{Y|X}(y|x) dy.$$

We now denote by  $m(X) = \mathbb{E}[g(Y)|X]$  the conditional expectation of  $g(Y)$  given  $X$ . This conditional expectation is a function of the rv  $X$  and therefore, is itself an rv, which takes the value  $m(x)$  when  $X = x$ . Hence

$$\begin{aligned} \mathbb{E}[\mathbb{E}[g(Y)|X]] &= E[m(X)] = \int m(x)f_X(x) dx \\ &= \int \left[ \int g(y)f_{Y|X}(y|x) dy \right] f_X(x) dx \\ &= \int g(y)f_Y(y) dy = E[g(Y)], \end{aligned}$$

because  $\int f_{Y|X}(y|x)fx(x) dx = \int f_{XY}(x,y) dx = f_Y(y)$  using Eq. (9). Next consider

$$\begin{aligned} E[h(X)g(Y)|X=x] &= \int h(x)g(y)f_{Y|X}(y|x) dy \\ &= h(x) \int g(y)f_{Y|X}(y|x) dy \\ &= h(x) E[g(Y)|X=x] \quad \text{for each } x. \end{aligned}$$

Hence  $E[h(X)g(Y)|X] = h(X)E[g(Y)|X]$ . We thus have the following important results:

$$EE[g(Y)|X] = E[g(Y)] \text{ and } E[h(X)g(Y)|X] = h(X)E[g(Y)|X].$$

We next consider the conditional properties of variance.

$$\begin{aligned} \text{Var}[Y] &= E[(Y - E[Y])^2] \\ &= EE[((Y - E[Y|X]) + (E[Y|X] - E[Y]))^2|X] \\ &= EE[(Y - E[Y|X])^2|X] + EE[(E[Y|X] - E[Y])^2|X] \\ &\quad + 2EE[(Y - E[Y|X])(E[Y|X] - E[Y])|X]. \end{aligned}$$

The three terms in the last expression are

$$EE[(Y - E[Y|X])^2|X] = E[\text{Var}[Y|X]],$$

since  $E[(Y - E[Y|X])^2|X] = \text{Var}[Y|X]$ ,

$$EE[(E[Y|X] - E[Y])^2|X] = E[(E[Y|X] - E[Y])^2] = \text{Var}[E[Y|X]],$$

and the third term is 0, using  $E[h(X)g(Y)|X] = h(X)E[g(Y)|X]$ .

*Summary.* Besides all the properties analogous to the properties of expectation, we have proved the following important properties of conditional expectation.

### Proposition 1.10.1.

- (i)  $E[E[g(Y)|X]] = E[g(Y)].$
- (ii)  $E[h(X)g(Y)|X] = h(X)E[g(Y)|X].$
- (iii)  $\text{Var}[Y] = E[\text{Var}[Y|X]] + \text{Var}[E[Y|X]].$

**Definition 1.10.1.** The function  $m(x) = E[Y|X=x]$  is called the regression function of  $Y$  on  $X$ . In particular, if  $m(x)$  is a linear function of  $x$ , then we can represent  $Y$  as

$$Y = \alpha + \beta X + \varepsilon \text{ with } E[\varepsilon|X] = 0,$$

and if  $\varepsilon$  is independent of  $X$ , then this is called the *linear regression model* and  $\text{Var}[\varepsilon]$ , if it exists, is the *residual variance*. More generally, if the dependence of  $Y$  on a  $k$ -dim rv  $X = (X_1, \dots, X_k)$  is such that  $m(x_1, \dots, x_k) = E[Y|X=x]$  is a linear function of  $(x_1, \dots, x_k)$ , then

$$Y = \alpha + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon \text{ with } E[\varepsilon|X] = 0,$$

and  $(X, Y)$  is said to follow a *multiple linear regression model* if  $\varepsilon$  is independent of  $X$ .

**Example 1.10.1.** The joint pdf of  $(X, Y)$  is

$$f_{XY}(x, y) = \begin{cases} C(x^2 + 2y^2) & 0 < x, y < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Find the constant  $C$ , the marginal pdf's of  $X$  and  $Y$ , the conditional pdf of  $Y$  given  $X = x$  and then find the means  $\mu_X$ ,  $\mu_Y$ , the variances  $\sigma_X^2$ ,  $\sigma_Y^2$ , the correlation coefficient  $\rho_{XY}$  and the conditional expectation  $E[Y|X = x]$ . Also find  $P[X > Y]$ .

*Solution.*

$$1 = C \int_0^1 \int_0^1 (x^2 + 2y^2) dx dy = C \left[ 1 \int_0^1 x^2 dx + 2 \int_0^1 y^2 dy \right] = C \left[ \frac{1}{3} + \frac{2}{3} \right].$$

Thus  $C = 1$  and  $f_{XY}(x, y) = x^2 + 2y^2$ ,  $0 < x < 1$ ,  $0 < y < 1$ . Now

$$\begin{aligned} f_X(x) &= \int_0^1 (x^2 + 2y^2) dy = x^2 + \frac{2}{3}, \quad 0 < x < 1, \\ f_Y(y) &= \int_0^1 (x^2 + 2y^2) dx = 2y^2 + \frac{1}{3}, \quad 0 < y < 1. \end{aligned}$$

The conditional pdf of  $Y$  given  $X = x$  is

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)} = \frac{x^2 + 2y^2}{x^2 + \frac{2}{3}}, \quad 0 < y < 1 \text{ for } 0 < x < 1.$$

Next, we evaluate the means, the variances, and the correlation coefficient:

$$\begin{aligned} \mu_X &= \int_0^1 x f_X(x) dx = \int_0^1 x \left( x^2 + \frac{2}{3} \right) dx \\ &= \frac{1}{4} + \frac{2}{3} \cdot \frac{1}{2} = \frac{7}{12}, \\ \mu_Y &= \int_0^1 y f_Y(y) dy = \int_0^1 y \left( 2y^2 + \frac{1}{3} \right) dy \\ &= 2 \cdot \frac{1}{4} + \frac{1}{3} \cdot \frac{1}{2} = \frac{2}{3}, \\ \sigma_X^2 &= E[X^2] - \mu_X^2 = \int_0^1 x^2 \left( x^2 + \frac{2}{3} \right) dx - \left( \frac{7}{12} \right)^2 \\ &= \left( \frac{1}{5} + \frac{2}{3} \cdot \frac{1}{3} \right) - \left( \frac{7}{12} \right)^2 = \frac{19}{45} - \frac{49}{144} = \frac{59}{720}, \\ \sigma_Y^2 &= E[Y^2] - \mu_Y^2 = \int_0^1 y^2 \left( 2y^2 + \frac{1}{3} \right) dy - \left( \frac{2}{3} \right)^2 \\ &= \left( \frac{2}{5} + \frac{1}{3} \cdot \frac{1}{3} \right) - \left( \frac{2}{3} \right)^2 = \frac{23}{45} - \frac{4}{9} = \frac{1}{15}, \\ \sigma_{XY} &= E[XY] - \mu_X \mu_Y = \int_0^1 \int_0^1 xy(x^2 + 2y^2) dx dy - \frac{7}{12} \cdot \frac{2}{3} \end{aligned}$$

$$= \left( \frac{1}{4} \cdot \frac{1}{2} + 2 \cdot \frac{1}{2} \cdot \frac{1}{4} \right) - \frac{7}{12} \cdot \frac{2}{3} = \frac{3}{8} - \frac{7}{18} = -\frac{1}{72},$$

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} = \frac{-\frac{1}{72}}{\sqrt{\frac{59}{720} \cdot \frac{1}{15}}} = -\frac{1}{72} \sqrt{\frac{10,800}{59}} = -0.1879.$$

The conditional expectation of  $Y$  given  $X = x$  is

$$\begin{aligned} E[Y|X = x] &= \int_0^1 y f_{Y|X}(y|x) dy = \int_0^1 y \frac{x^2 + 2y^2}{x^2 + \frac{2}{3}} dy \\ &= \frac{1}{x^2 + \frac{2}{3}} \int_0^1 (x^2 y + 2y^3) dy = \frac{x^2 \frac{1}{2} + 2 \cdot \frac{1}{4}}{x^2 + \frac{2}{3}} \\ &= \frac{x^2 + 1}{2x^2 + \frac{4}{3}}. \end{aligned}$$

Note that

$$\begin{aligned} E[E[Y|X]] &= \int_0^1 E[Y|X = x] f_X(x) dx = \int_0^1 \frac{x^2 + 1}{2x^2 + \frac{4}{3}} \left( x^2 + \frac{2}{3} \right) dx \\ &= \int_0^1 \frac{1}{2} (x^2 + 1) dx = \frac{1}{2} \left( \frac{1}{3} + 1 \right) = \frac{2}{3} = \mu_Y, \text{ as it should be.} \end{aligned}$$

Finally,

$$P[X > Y] = \int_0^1 \left[ \int_0^x (x^2 + 2y^2) dy \right] dx = \int_0^1 \left( x^2 + 2 \frac{x^3}{3} \right) dx = \frac{5}{3} \cdot \frac{1}{4} = \frac{5}{12}.$$

## 1.11 Transforms of Random Variables and Their Distributions

We start with some simple transforms. In our following discussion, we shall often work with continuous rv's and their pdf's, of which (i), (ii), and (iii) will remain valid for discrete rv's and their pmf's by replacing  $\int dx$  by  $\sum_x$ .

- (i)** For  $b > 0$ ,  $F_{a+bX}(z) = P[a + bX \leq z] = P[X \leq \frac{z-a}{b}] = F_X(\frac{z-a}{b})$ . Hence  $f_{a+bX}(z) = b^{-1} f_X(\frac{z-a}{b})$ . More generally,

$$\text{for } b \neq 0, f_{a+bX}(z) = |b|^{-1} f_X\left(\frac{z-a}{b}\right). \quad (10)$$

In the discrete case,  $f_{a+bX}(z) = f_X(\frac{z-a}{b})$ .

- (ii)**  $F_{X^2}(z) = P[-\sqrt{z} \leq X \leq \sqrt{z}] = F_X(\sqrt{z}) - F_X(-\sqrt{z})$  for  $z > 0$ . Hence

$$\begin{aligned} f_{X^2}(z) &= \frac{1}{2\sqrt{z}} f_X(\sqrt{z}) - \left( -\frac{1}{2\sqrt{z}} \right) f_X(-\sqrt{z}) \\ &= \frac{1}{2\sqrt{z}} [f_X(\sqrt{z}) + f_X(-\sqrt{z})] \text{ for } z > 0. \end{aligned} \quad (11a)$$

If  $X$  is distributed symmetrically about 0, then  $f_X(\sqrt{z}) = f_X(-\sqrt{z})$ , so that

$$f_{X^2}(z) = \frac{1}{\sqrt{z}} f_X(\sqrt{z}), \quad z > 0 \quad (11b)$$

for symmetrically distributed  $X$ . In the discrete case,

$$f_{X^2}(z) = f_X(\sqrt{z}) + f_X(-\sqrt{z})$$

and  $f_{X^2}(z) = 2f_X(\sqrt{z})$  if  $X$  is symmetric.

(iii) If  $X_1, X_2$  are independent, then

$$\begin{aligned} F_{X_1+X_2}(z) &= \int \int_{x_1+x_2 \leq z} f_{X_1}(x_1)f_{X_2}(x_2) dx_1 dx_2 \\ &= \int_{x_1=-\infty}^{\infty} \int_{x_2=-\infty}^{z-x_1} f_{X_1}(x_1)f_{X_2}(x_2) dx_1 dx_2. \end{aligned}$$

Differentiating under the integral, we have

$$\begin{aligned} f_{X_1+X_2}(z) &= \int_{-\infty}^{\infty} f_{X_1}(x_1) \left[ \frac{d}{dz} \int_{-\infty}^{z-x_1} f_{X_2}(x_2) dx_2 \right] dx_1 \\ &= \int_{-\infty}^{\infty} f_{X_1}(x) f_{X_2}(z-x) dx. \end{aligned}$$

If  $X_1, X_2$  are nonnegative, then  $f_{X_1}(x) = 0$  for  $x < 0$  and  $f_{X_2}(z-x) = 0$  for  $x > z$ . In this case the above formula becomes

$$f_{X_1+X_2}(z) = \int_0^z f_{X_1}(x) f_{X_2}(z-x) dx. \quad (12a)$$

For independent discrete rv's,

$$f_{X_1+X_2}(z) = \sum_x f_{X_1}(x) f_{X_2}(z-x). \quad (12b)$$

(iv) If  $X_1, X_2$  are independent with  $P[X_2 > 0] = 1$ , then

$$F_{\frac{X_1}{X_2}}(z) = \int_{x_2=0}^{\infty} \left[ \int_{x_1=-\infty}^{x_2 z} f_{X_1}(x_1) dx_1 \right] f_{X_2}(x_2) dx_2$$

and differentiating under the integral, we have

$$f_{\frac{X_1}{X_2}}(z) = \int_0^{\infty} x f_{X_1}(xz) f_{X_2}(x) dx. \quad (13)$$

(v) For mutually independent  $X_1, \dots, X_k$ , let  $Y = \min(X_1, \dots, X_k)$  and  $Z = \max(X_1, \dots, X_k)$ . Then

$$\begin{aligned} 1 - F_Y(y) &= P[\min(X_1, \dots, X_k) > y] \\ &= \prod_{i=1}^k P[X_i > y] = \prod_{i=1}^k \{1 - F_{X_i}(y)\} \text{ and} \\ F_Z(z) &= P[\max(X_1, \dots, X_k) \leq z] \\ &= \prod_{i=1}^k P[X_i \leq z] = \prod_{i=1}^k F_{X_i}(z). \end{aligned}$$

Hence

$$\begin{aligned} f_Y(y) &= -\frac{d}{dy} \prod_{i=1}^k \{1 - F_{X_i}(y)\} = \sum_{i=1}^k f_{X_i}(y) \prod_{j \neq i} [1 - F_{X_j}(y)] \\ f_Z(z) &= \frac{d}{dz} \prod_{i=1}^k F_{X_i}(z) = \sum_{i=1}^k f_{X_i}(z) \prod_{j \neq i=1}^k F_{X_j}(z). \end{aligned}$$

We now consider a transformation of

$$\mathbf{X} = (X_1, \dots, X_k) \xrightarrow{g} (Y_1, \dots, Y_k) = \mathbf{Y}$$

(ie,  $Y_i = g_i(X_1, \dots, X_k)$ ), where  $g$  is one to one with continuous partial derivatives and a nonvanishing Jacobian of transformation in the continuous case. Let  $g^{-1}$  denote the inverse transformation, that is,

$$g^{-1}(\mathbf{y}) = g^{-1}(y_1, \dots, y_k) := (h_1(y_1, \dots, y_k), \dots, h_k(y_1, \dots, y_k)).$$

In the discrete case,

$$f_{\mathbf{Y}}(y_1, \dots, y_k) = f_{\mathbf{X}}(h_1(y_1, \dots, y_k), \dots, h_k(y_1, \dots, y_k)) = f_{\mathbf{X}}(g^{-1}(\mathbf{y})).$$

In the continuous case, for all events  $B$ ,  $P[Y \in B] = P[\mathbf{X} \in g^{-1}(B)]$ , where  $g^{-1}(B) = \{\mathbf{x}: g(\mathbf{x}) \in B\}$ . Thus

$$\begin{aligned} \int_B f_{\mathbf{Y}}(y_1, \dots, y_k) dy_1 \cdots dy_k \\ &= \int_{g^{-1}(B)} f_{\mathbf{X}}(x_1, \dots, x_k) dx_1 \cdots dx_k \\ &= \int_{g(g^{-1}(B))} f_{\mathbf{X}}(g^{-1}(\mathbf{y})) \left| \det \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \cdots & \frac{\partial x_k}{\partial y_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_1}{\partial y_k} & \cdots & \frac{\partial x_k}{\partial y_k} \end{bmatrix} \right| dy_1 \cdots dy_k \\ &= \int_B f_{\mathbf{X}}(g^{-1}(\mathbf{y})) |J_{g^{-1}(\mathbf{y})}| dy_1 \cdots dy_k, \end{aligned}$$

where  $J_{g^{-1}(\mathbf{y})} = \det \begin{bmatrix} \frac{\partial h_1}{\partial y_1} & \cdots & \frac{\partial h_k}{\partial y_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_1}{\partial y_k} & \cdots & \frac{\partial h_k}{\partial y_k} \end{bmatrix}$ , and note that  $J_{g^{-1}(\mathbf{y})} = [J_g(g^{-1}(\mathbf{y}))]^{-1}$ , so we compute the one which is easier. Thus

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(g^{-1}(\mathbf{y})) |J_{g^{-1}(\mathbf{y})}| = f_{\mathbf{X}}(g^{-1}(\mathbf{y})) \frac{1}{|J_g(g^{-1}(\mathbf{y}))|}.$$

## An Extension

Let  $S_1, \dots, S_l$  be disjoint open sets in  $\mathbb{R}^k$  with  $\sum_{j=1}^l P[X \in S_j] = 1$ . Let  $g: \bigcup_{j=1}^l S_j \rightarrow \mathbb{R}^k$ ,  $\mathbf{Y} = g(\mathbf{X})$ , where for each  $j$ , the restriction  $g_j$  of  $g$  on  $S_j$  is one to one with continuous partial derivatives and nonvanishing  $J_{g_j}$ . Then for all events  $B$ ,

$$\begin{aligned} \int_B f_Y(\mathbf{y}) d\mathbf{y} &= P[\mathbf{Y} \in B] = \sum_{j=1}^l P[\mathbf{Y} \in B, \mathbf{X} \in S_j] \\ &= \sum_{j=1}^l P[X \in g_j^{-1}(B), \mathbf{X} \in S_j] \\ &= \sum_{j=1}^l \int_{g_j^{-1}(B) \cap S_j} f_X(\mathbf{x}) d\mathbf{x} \\ &= \sum_{j=1}^l \int_{B \cap g_j(S_j)} f_X(g^{-1}(\mathbf{y})) |J_{g^{-1}}(\mathbf{y})| d\mathbf{y} \\ &= \int_B \sum_{j=1}^l f_X(g^{-1}(\mathbf{y})) |J_{g^{-1}}(\mathbf{y})| I_{g(S_j)}(\mathbf{y}) d\mathbf{y}, \end{aligned}$$

where  $I_A(\mathbf{y}) = 1$  if  $\mathbf{y} \in A$  and  $= 0$  if  $\mathbf{y} \notin A$ . Hence

$$f_Y(\mathbf{y}) = \sum_{j=1}^l f_X(g^{-1}(\mathbf{y})) |J_{g^{-1}}(\mathbf{y})| I_{g(S_j)}(\mathbf{y}).$$

## Joint Distribution of Order Statistics

Let  $X_1, \dots, X_n$  be iid real-valued rv's with common pdf  $f$ . Define

$$Y_i = X_{n:i}, \quad i = 1, \dots, n,$$

where  $X_{n:1} < \dots < X_{n:n}$  are the order statistics, that is, the ordered values of  $X_1, \dots, X_n$ . (The  $X'_i$ 's are all distinct with probability 1.) The joint pdf of  $(X_1, \dots, X_n)$  is  $f_X(\mathbf{x}) = \prod_{i=1}^n f(x_i)$ . For the  $n!$  permutations  $(j_1, \dots, j_n)$  of  $(1, \dots, n)$ , let

$$S_{j_1 \dots j_n} = \{(x_1, \dots, x_n) : x_{j_1} < \dots < x_{j_n}\}$$

and define  $\mathbf{Y} = (Y_1, \dots, Y_n) = g_{j_1, \dots, j_n}(X_1, \dots, X_n) = (X_{j_1}, \dots, X_{j_n})$  on  $S_{j_1, \dots, j_n}$  for each  $(j_1, \dots, j_n)$ . Then  $Y_1 < \dots < Y_n$  with probability 1 and  $g_{j_1, \dots, j_n}$  is one to one with  $|J_{g_{j_1, \dots, j_n}}(\mathbf{y})| = 1$  for each permutation  $(j_1, \dots, j_n)$  of  $(1, \dots, n)$ . Hence

$$f_Y(\mathbf{y}) = \sum_{j_1, \dots, j_n} f_X(g_{j_1, \dots, j_n}^{-1}(\mathbf{y})) = n! \prod_{i=1}^n f(y_i) \quad \text{for } y_1 < \dots < y_n$$

and  $f_Y(\mathbf{y}) = 0$  otherwise.

## Linear Transformation

As a special case of  $g: \mathbf{X} \rightarrow \mathbf{Y}$  from  $\mathbb{R}^k \rightarrow \mathbb{R}^k$ , consider a linear transformation  $\mathbf{Y} = \mathbf{AX}$  where  $\mathbf{A}$  is a nonsingular  $k \times k$  matrix. Then the inverse transformation is  $\mathbf{X} = \mathbf{A}^{-1}\mathbf{Y} = g^{-1}(\mathbf{Y})$ , so  $|J_{g^{-1}}(\mathbf{y})| = |\det(\mathbf{A}^{-1})| = \frac{1}{|\det(\mathbf{A})|}$ . Hence

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(\mathbf{A}^{-1}\mathbf{y}) \frac{1}{|\det(\mathbf{A})|}. \quad (14)$$

In particular, if  $\mathbf{A}$  is orthonormal, that is,  $\mathbf{AA}^T = \mathbf{A}^T\mathbf{A} = \mathbf{I}$ , then  $|\det(\mathbf{A})| = 1$ , so that  $f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(\mathbf{A}^{-1}\mathbf{y})$ .

## Probability Integral Transform

Suppose  $X$  has cdf  $F$  which is continuous and strictly increasing. Then  $F^{-1}$  is uniquely defined as

$$F^{-1}(u) = x \text{ iff } F(x) = u \text{ for } 0 < u < 1.$$

Then the cdf of  $Y = F(X)$  at  $u \in (0, 1)$  is

$$F_Y(u) = P[F(x) \leq u] = P[X \leq F^{-1}(u)] = F(F^{-1}(u)) = u.$$

Thus  $f_Y(u) = 1$  for  $0 < u < 1$  and  $f_Y(u) = 0$  for  $u \notin (0, 1)$ , because  $0 < Y = F(X) < 1$  with probability 1. In other words, if  $X$  has a continuous and strictly increasing cdf  $F$ , then  $Y = F(X)$  is distributed with pdf

$$f_Y(u) = \begin{cases} 1 & \text{if } 0 < u < 1, \\ 0 & \text{otherwise.} \end{cases}$$

A rv with this pdf is said to be a Uniform(0, 1) rv. Conversely, if  $U$  is Uniform(0, 1), then  $X = F^{-1}(U)$  has cdf  $F$ . This fact is useful in generating random samples (ie, iid rv's) with cdf  $F$  by first generating random samples  $U_1, U_2, \dots$  from Uniform(0, 1), which is easy, and then transforming  $U_1, U_2, \dots$  to  $X_1 = F^{-1}(U_1), X_2 = F^{-1}(U_2), \dots$

## Exercises

- 1.1. Prove that the axioms P(i, ii), P(i, iii), and P(i, iv) are equivalent.
- 1.2. Prove [Proposition 1.4.2](#) by induction, starting from [Proposition 1.4.1](#)(iii).
- 1.3. Five cards with numbers 1, ..., 5 are arranged in random order. A person claiming psychic power, declares the arrangement. Assuming that the numbers are declared purely at random, what is the probability that
  - (a) All are correct? (b) Exactly 4 are correct?
- 1.4. From a well-shuffled standard deck of cards, eight cards are dealt.
  - (a) What is the probability that the hand contains two trios (exactly three cards of the same denomination such as three kings)?

- (b) What is the probability that the hand contains a run of five cards (exactly five cards of consecutive denominations such as 8, 9, 10, J, Q irrespective of suits, 10, J, Q, K, A included)?
- 1.5. A poker hand consists of five cards dealt from a well-shuffled standard deck of cards. A hand is a
- (a) *straight* if the cards are of consecutive denominations (including A2345 and AKQJ10), but not of the same suit,
  - (b) *flush* if all five cards are of the same suit but not a straight,
  - (c) *full house* if three cards are of the same denomination and the other two cards are also of the same denomination such as three 10s and two Queens.
- Find the probabilities of a hand being a straight, a flush, and a full house.
- 1.6. (*The birthday problem*) In an assembly of  $n$  people, what is the probability that no two have the same birthday? Find the smallest  $n$  for which this probability is less than 1/2. [For simplicity, ignore leap-years and assume that the birthdays are equally likely to be any of the 365 days of a year.]
- 1.7. Show that for any event  $B$  with  $P[B] > 0$ ,  $P[A|B]$  as a function of  $A \in \mathcal{A}$  is a probability.
- 1.8. Prove [Proposition 1.6.2](#).
- 1.9. A couple has two children. If one child is a boy, what is the probability that the other child is a girl?
- 1.10. Two shooters  $A$  and  $B$  can hit a target with probabilities 0.8 and 0.7, respectively. They shoot alternatively and the one who first makes three hits wins. Find the probability that  $A$  wins if
- (a)  $A$  shoots first, (b)  $B$  shoots first, (c) the one who shoots first is decided by the toss of a fair coin.
- 1.11. There are two dice in a bag both of which are balanced, but one has 1, 2, 3, 4, 5, 6 and the other has 1, 1, 2, 3, 4, 5 on the six faces. One of the dice is selected at random and rolled.
- (a) What is the probability that the outcome is an even number?
  - (b) If the outcome is an even number, what is the probability that the standard die was rolled?
- 1.12. In a city, in the month of July, the maximum temperature reaches  $100^{\circ}\text{F}$  or above with probability 0.2 on any day and with probability 0.8 if the previous day was  $100^{\circ}\text{F}$  or above. What is the probability that there will be exactly 5 consecutive days of  $100^{\circ}\text{F}$  or above starting on July 15?
- 1.13. Prove [Proposition 1.7.1](#).
- 1.14. There are  $N$  tickets in a bag, numbered  $1, \dots, N$ , from which  $n$  tickets are drawn. Let  $X$  be the largest number drawn. Find the pmf and the cdf of  $X$  if the tickets are drawn at random (a) with replacement, (b) without replacement.
- 1.15. The lifetime of an equipment (in hours) is a random variable  $X$  with pdf

$$f(x) = \begin{cases} c & 0 < x \leq 50 \\ c \exp(-(x - 50)) & x > 50. \end{cases}$$

Find (a) the constant  $c$ , (b) the cdf of  $X$ , (c)  $P[10 < X \leq 100]$ .

- 1.16.** Let  $X$  be a random variable with pmf

$$f(x) = c/2^x, \quad x = 1, 2, \dots$$

Find (a) the constant  $c$ , (b) the cdf of  $X$ , (c)  $P[X \geq 5]$ .

- 1.17.** A game consists of tossing a fair coin and then drawing a random number from  $0, 1, \dots, 9$  with equal probabilities if head comes up, or spinning a wheel to choose a random number in the interval  $(0, 10]$  if tail comes up. Let  $X$  be the outcome of this game, which is either an integer between 0 and 9 or a number in the interval  $(0, 10]$ .

- (a) Find the cdf of  $X$  and  $P[3 \leq X \leq 5]$ . Note that the rv  $X$  is neither discrete nor continuous, but a mixture.
- (b) Express  $X$  as a mixture,  $X = \alpha U + (1 - \alpha)V$ , where  $U$  is discrete and  $V$  is continuous, giving the pmf of  $U$  and the pdf of  $V$ . [Interpretation of mixture is  $F_X = \alpha F_U + (1 - \alpha)F_V$ .]

- 1.18.** (*Censoring: another example of a mixture*) Let  $T$  be a nonnegative rv (life of an equipment or survival time of a patient undergoing a certain treatment) and  $t_c > 0$  is a time at which observation stops. So we observe  $T$  if  $0 \leq T < t_c$  and  $t_c$  if  $T \geq t_c$ . The resulting observation is  $Y = \min(T, t_c)$ . Find the cdf of  $Y$  which is a mixture.

- 1.19.** Verify (i)–(vi) of [Proposition 1.8.1](#).

- 1.20.** Verify (i)–(vi) of [Proposition 1.8.2](#).

- 1.21.** Let  $(X, Y)$  be a two-dimensional rv with joint pdf

$$f(x, y) = c(x + 2y), \quad 0 \leq x, \quad y \leq 1.$$

- (a) Find the constant  $c$  and  $P[X \leq Y]$ .

- (b) Find the marginal pdf's  $f_X$  of  $X$  and  $f_Y$  of  $Y$ .

- (c) Calculate  $\mu_X, \mu_Y, \sigma_X^2$ , and  $\sigma_Y^2$ . Also find the medians of  $X$  and  $Y$ .

- (d) Find  $\text{Cov}[X, Y]$  and  $\rho_{XY}$ .

- 1.22.** The joint pdf of  $(X, Y)$  is

$$f(x, y) = c(2x^2 + xy), \quad 0 \leq x \leq 1, \quad 0 \leq y \leq 2.$$

Find

- (a) the constant  $c$ , (b) the marginal pdf  $f_X$  of  $X$ , (c) the conditional pdf  $f_{Y|X}(y|x)$  of  $Y$  given  $X = x$ , (d) the regression function  $m(x) = E[Y|X = x]$ .

- 1.23.** The pdf of  $X$  and the conditional pdf of  $Y$  given  $X = x$  are given by

$$f_X(x) = \exp(-x), \quad 0 < x < \infty, \quad \text{and } f_{Y|X}(y|x) = x \exp(-xy), \quad 0 < y < \infty.$$

- (a) Find the joint pdf of  $(X, Y)$ , the marginal pdf of  $Y$ , and the conditional pdf of  $X$  given  $Y = y$ .

- (b) Find the regression function  $m(x) = E[Y|X = x]$ .

**1.24.** The joint pdf of  $(X, Y)$  is

$$f(x, y) = cx \exp(-x), \quad 0 < y < x < \infty.$$

Find

- (a)** the constant  $c$  and the marginal pdf's  $f_X$  and  $f_Y$ ,
- (b)**  $E[X]$ ,  $\text{Var}[X]$ ,  $E[Y]$ ,  $\text{Var}[Y]$ ,  $\text{Cov}[X, Y]$ , and  $\rho_{XY}$ ,
- (c)** the conditional pdf of  $Y$  given  $X = x$ ,  $E[Y|X = x]$  and  $\text{Var}[Y|X = x]$ , and
- (d)** verify:  $E[Y] = E[E(Y|X)]$  and  $\text{Var}[Y] = E[\text{Var}(Y|X)] + \text{Var}[E(Y|X)]$ .

# Some Common Probability Distributions

## 2.1 Discrete Distributions

We shall start with the simplest nontrivial rv  $X$  which can take only one of two values, say 0 and 1, with  $P[X = 0] = 1 - p = q$  and  $P[X = 1] = p$  where  $0 < p < 1$ . This arises in the context of a random experiment resulting in success or failure, such as the toss of a coin resulting in head (success) or tail (failure), or a medicine having favorable effect (success) or not (failure) on a patient, etc., with  $P[\text{success}] = p$  and  $P[\text{failure}] = q = 1 - p$ . The rv  $X$  takes the values 1 or 0 for the outcomes success or failure, respectively. Such an rv is called a *Bernoulli*( $p$ ) rv and we write the pmf of  $X$  as

$$f_X(x) = p^x(1-p)^{1-x}, \quad x = 0, 1. \quad (1)$$

### 2.1.1 Binomial Distribution $\text{Bin}(n, p)$

Suppose  $X_1, \dots, X_n$  are independent *Bernoulli*( $p$ ) rv's as described by Eq. (1) and let  $X = X_1 + \dots + X_n$ , that is,  $X$  = total number of successes in  $n$  independent experiments with  $P[\text{success}] = p$  in each experiment. The pmf of  $X$  is obtained by noting that each sequence of  $n$  outcomes with  $x$  successes and  $n - x$  failures has probability  $= p^x(1-p)^{n-x}$ , and the number of such sequences is the number of ways in which  $x$  trials (resulting in successes) can be chosen from  $n$  trials, which is  $\binom{n}{x} = \frac{n!}{x!(n-x)!}$ . Hence the pmf of  $X$  is

$$f_X(x) = \binom{n}{x} p^x(1-p)^{n-x}, \quad x = 0, 1, \dots, n.$$

Since  $E[X_1] = 1 \cdot p + 0 \cdot (1 - p) = p$  and

$$\begin{aligned} \text{Var}[X_1] &= E[X_1^2] - (E[X_1])^2 \\ &= [1^2 \cdot p + 0^2 \cdot (1 - p)] - p^2 \\ &= p - p^2 = p(1 - p), \end{aligned}$$

we have

$$E[X] = nE[X_1] = np \text{ and } \text{Var}[X] = n \text{Var}[X_1] = np(1 - p) \quad (2)$$

**Example 2.1.1.** An equipment with  $n$  components needs at least 3 to function properly for its overall effectiveness. If each of these components function with probability  $4/5$ , independently of one another, what should be the minimum  $n$  for the equipment to remain effective with probability 0.9 or more?

*Solution.*

$$\begin{aligned} P[\text{Not Effective}] &= \sum_{x=0}^2 \binom{n}{x} \left(\frac{4}{5}\right)^x \left(\frac{1}{5}\right)^{n-x} \\ &= \frac{1}{5^n} \left[ \binom{n}{0} 4^0 + \binom{n}{1} 4^1 + \binom{n}{2} 4^2 \right] \\ &= \frac{1}{5^n} [1 + 4n + 8n(n-1)] \\ &= \frac{1}{5^n} [1 - 4n + 8n^2]. \end{aligned}$$

We want the smallest  $n$  so that  $\frac{1}{5^n} [1 - 4n + 8n^2] \leq \frac{1}{10}$ .

$n$	$(1 - 4n + 8n^2)/5^n$
3	$(1 - 12 + 72)/5^3 = 0.488$
4	$(1 - 16 + 128)/5^4 = 0.1808$
5	$(1 - 20 + 200)/5^5 = 0.05792$

so the minimum  $n$  to keep the equipment effective,  $n_{\min} = 5$ .

### 2.1.2 Multinomial Distribution $Multi(n; p_1, \dots, p_k)$

Extending the concept of the  $Binomial(n, p)$  distribution from the outcomes of  $n$  independent trials, each with two possible outcomes, consider  $n$  independent trials with  $k$  possible outcomes  $1, \dots, k$  with

$$P[i\text{th outcome}] = p_i > 0, \quad \sum_{i=1}^k p_i = 1.$$

Let  $X = (X_1, \dots, X_k)$  with  $X_i$  = number of trials resulting in the  $i$ th outcome. Then any sequence of  $n$  outcomes with  $x_i$  = number of trials with the  $i$ th outcome,  $i = 1, \dots, k$  has probability

$$p_1^{x_1} \cdots p_k^{x_k} \text{ if } x_i \geq 0, \quad \sum_{i=1}^k x_i = n \text{ and } 0 \text{ otherwise,}$$

and the number of such sequences is the number of ways in which such  $(x_1, \dots, x_k)$  can be chosen from the  $n$  trials. This number is  $n!/(x_1! \cdots x_k!)$ , the multinomial coefficient. Thus, the joint pmf of  $X = (X_1, \dots, X_k)$  is

$$f_X(x_1, \dots, x_k) = \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \cdots p_k^{x_k}$$

if  $x_1, \dots, x_k \geq 0$  and  $\sum_{i=1}^k x_i = n$ , and 0 otherwise.

The multinomial coefficient is obtained by multiplying the number of successive choices of  $x_1$  out of  $n$  trials, followed by  $x_2$  out of  $n - x_1$  trials and so on, and finally,  $x_{k-1}$  out of  $n - (x_1 + \dots + x_{k-2})$ . Thus, using  $x_k = n - \sum_{i=1}^{k-1} x_i$ , we have

$$\begin{aligned} & \binom{n}{x_1} \binom{n-x_1}{x_2} \cdots \binom{n-(x_1+\dots+x_{k-2})}{x_{k-1}} \\ &= \frac{n!}{x_1!(n-x_1)!} \cdot \frac{(n-x_1)!}{x_2!(n-x_1-x_2)!} \cdots \frac{(n-x_1-\dots-x_{k-2})!}{x_{k-1}!x_k!}. \end{aligned}$$

By cancellation, this yields the multinomial coefficient. To find  $E[X_i]$ ,  $\text{Var}[X_i]$ , and  $\text{Cov}[X_i, X_j]$  for any  $i \neq j$ , note that each  $X_i \sim \text{Bin}(n, p_i)$  and  $X_j | X_i \sim \text{Bin}(n - X_i, p_j / (1 - p_i))$ . Hence

$$E[X_i] = np_i, \quad \text{Var}[X_i] = np_i(1 - p_i) \text{ and}$$

as in Eq. (2), and using Proposition 1.10.1,

$$\begin{aligned} \text{Cov}[X_i, X_j] &= E[X_i X_j] - E[X_i]E[X_j] \\ &= E[X_i E(X_j | X_i)] - n^2 p_i p_j \\ &= E[X_i (n - X_i) p_j / (1 - p_i)] - n^2 p_i p_j \\ &= \frac{p_j}{1 - p_i} [n \cdot np_i - \{np_i(1 - p_i) + n^2 p_i^2\}] - n^2 p_i p_j \\ &= -np_i p_j \end{aligned}$$

by algebraic simplification.

### 2.1.3 Geometric Distribution $\text{Geom}(p)$

In a sequence of independent trials, each resulting in success with probability  $p$  or failure with probability  $1 - p = q$ , let  $X$  be the number of trials needed for the first success to occur. Then the pmf of  $X$  is

$$\begin{aligned} f_X(n) &= P[\text{Failures in the first } n-1 \text{ trials and success in the } n\text{th trial}] \\ &= q^{n-1} p, \quad n = 1, 2, \dots \end{aligned}$$

Then  $\sum_{n=1}^{\infty} f_X(n) = \sum_{n=1}^{\infty} q^{n-1} p = p \sum_{s=0}^{\infty} q^s = p(1 - q)^{-1} = 1$ , and

$$\begin{aligned}
E[X] &= \sum_{n=1}^{\infty} nf_X(n) = \sum_{n=1}^{\infty} nq^{n-1}p = p \frac{d}{dq} \sum_{s=0}^{\infty} q^s \\
&= p \cdot \frac{d}{dq} (1-q)^{-1} = p(1-q)^{-2} = \frac{1}{p} \\
E[X^2] &= \sum_{n=1}^{\infty} n^2 f_X(n) = \sum_{n=1}^{\infty} n^2 q^{n-1}p = p \sum_{n=1}^{\infty} [n + n(n-1)] q^{n-1} \\
&= p \sum_{n=1}^{\infty} nq^{n-1} + pq \sum_{n=1}^{\infty} n(n-1)q^{n-2} = E[X] + pq \frac{d^2}{dq^2} \sum_{s=0}^{\infty} q^s \\
&= \frac{1}{p} + pq \frac{d^2}{dq^2} (1-q)^{-1} = \frac{1}{p} + pq \cdot 2(1-q)^{-3} = \frac{2}{p^2} - \frac{1}{p}. \\
\text{Var}[X] &= E[X^2] - (E[X])^2 = \left( \frac{2}{p^2} - \frac{1}{p} \right) - \frac{1}{p^2} = \frac{q}{p^2}.
\end{aligned}$$

*Remark 2.1.1.* Memoryless property of the Geometric distribution. If  $X \sim \text{Geom}(p)$ , then  $P[X > k] = p \sum_{n=k+1}^{\infty} q^{n-1} = pq^k \sum_{s=0}^{\infty} q^s = pq^k (1-q)^{-1} = q^k$ . More directly,  $P[X > k] = P[\text{all failures in the first } k \text{ trials}] = q^k$ . So

$$P[X = k+n | X > k] = \frac{P[X = k+n]}{P[X > k]} = \frac{q^{k+n-1}p}{q^k} = q^{n-1}p = P[X = n].$$

## 2.1.4 Negative Binomial Distribution $N\text{Bin}(r, p)$

In a sequence of independent trials resulting in success or failure with probabilities  $p$  and  $q = 1 - p$  respectively, let  $X$  be the number of trials needed for the  $r$ th success to occur. Then

$$\begin{aligned}
f_X(n) &= P[(r-1) \text{ successes in the first } (n-1) \text{ trials and success on the } n\text{th trial}] \\
&= \binom{n-1}{r-1} p^{r-1} q^{n-r}, \quad n = r, r+1, \dots
\end{aligned}$$

Note that this  $X$  can be expressed as  $X = Y_1 + \dots + Y_r$ , where  $Y_1, \dots, Y_r$  are independent  $\text{Geom}(p)$ . Hence

$$E[X] = rE[Y_1] = r/p \text{ and } \text{Var}[X] = r \text{ Var}[Y_1] = rq/p^2.$$

**Note.** A binomial rv is the number of successes in a given number of trials, whereas, a negative binomial rv is the number of trials needed for a given number of successes.

**Example 2.1.2.** A target-shooter can hit the bull's eye once in three attempts on average, that is, with probability  $1/3$ . In a tournament with an entry fee of \$10, there are prizes of \$100, \$50, and \$20 if one can hit the bull's eye three times in 3, 4, or 5 attempts, respectively. Should our target-shooter enter this tournament paying the fee?

*Solution.* Let  $N$  = number of attempts needed to hit the bull's eyes  $r = 3$  times, where the probability of a hit is  $p = 1/3$  and the results of the attempts are mutually independent,

and for  $N = n$ , let  $X = X(n)$  be the prize won. Then  $X(3) = \$100$ ,  $X(4) = \$50$ ,  $X(5) = \$20$ ,  $X(n) = 0$  for  $n > 5$ , and  $N$  is  $NB(r = 3, p = 1/3)$ . Since

$$P[N = n] = \binom{n-1}{r-1} p^r (1-p)^{n-r} = \binom{n-1}{2} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^{n-3},$$

$$P[N = 3] = \binom{2}{2} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^0 = \frac{1}{27}$$

$$P[N = 4] = \binom{3}{2} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^1 = \frac{2}{27}$$

$$P[N = 5] = \binom{4}{2} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^2 = \frac{8}{81}.$$

So  $E[X] = \$100 \cdot \frac{1}{27} + \$50 \cdot \frac{2}{27} + \$20 \cdot \frac{8}{81} = \$\frac{760}{81} \approx \$9.38$ . Since the entry fee of \$10 is greater than the expected value of prize to be won, which is \$9.38, this game is not in favor of our target-shooter, who (strictly speaking) should not enter unless for fun, because the expected loss is very little.

### 2.1.5 Hypergeometric Distribution ( $n, N, m$ )

Suppose that from a box containing  $N$  balls of which  $m$  are white and  $N - m$  are red,  $n$  balls are drawn at random *without replacement*. Then the number of white balls in the sample is an rv  $X$  with pmf

$$f_X(x) = \binom{m}{x} \binom{N-m}{n-x} / \binom{N}{n}, \quad x = 0, \dots, n,$$

with the convention  $\binom{r}{k} = 0$  for  $k < 0$  or  $k > r$ . Thus the effective range of  $X$  is  $\max(0, m+n-N) \leq x \leq \min(m, n)$ . To find the mean and variance of  $X$ , we calculate  $E[X^k]$  for  $k = 1$  and  $k = 2$ , using  $x \binom{m}{x} = m \binom{m-1}{x-1}$  and  $n \binom{N}{n} = N \binom{N-1}{n-1}$  as shown below.

$$\begin{aligned} E[X^k] &= \sum_{x=0}^n x^k \cdot \frac{\binom{m}{x} \binom{N-m}{n-x}}{\binom{N}{n}} \\ &= \sum_{x=1}^n x^{k-1} \cdot \frac{m \binom{m-1}{x-1} \binom{N-m}{n-x}}{\frac{N}{n} \binom{N-1}{n-1}} \\ &= \frac{mn}{N} \sum_{x-1=0}^{n-1} x^{k-1} \cdot \frac{\binom{m-1}{x-1} \binom{N-m}{n-x}}{\binom{N-1}{n-1}} \end{aligned}$$

$$\begin{aligned}
&= \frac{mn}{N} \sum_{y=0}^{n-1} (y+1)^{k-1} \cdot \frac{\binom{m-1}{y} \binom{N-m}{n-1-y}}{\binom{N-1}{n-1}} \\
&= \frac{mn}{N} E[(Y+1)^{k-1}]
\end{aligned}$$

where  $Y$  is Hypergeometric  $(n-1, N-1, m-1)$ . Hence

$$E[X] = \frac{mn}{N} E[(Y+1)^0] = \frac{mn}{N} \cdot 1 = \frac{mn}{N} = np,$$

where  $p = m/N$  = proportion of white balls in the box, and

$$\begin{aligned}
\text{Var}[X] &= E[X^2] - (E[X])^2 \\
&= np\{E[Y] + 1\} - (np)^2 \\
&= np \left\{ \frac{(m-1)(n-1)}{N-1} + 1 \right\} - (np)^2 \\
&= \frac{N-n}{N-1} np(1-p)
\end{aligned}$$

by algebraic simplification.

*Remark 2.1.2.*

1. If in the above set-up,  $n$  balls are drawn at random with replacement , then the number of white balls in the sample is distributed as  $\text{Bin}(n, p)$ . In the without replacement case, the mean of the hypergeometric rv remains the same as in  $\text{Bin}(n, p)$ , but the variance is reduced by a factor of  $\frac{N-n}{N-1}$ .
2. For any  $n$ , the factor  $\frac{N-n}{N-1} \rightarrow 1$  as  $N \rightarrow \infty$ , which agrees with the common sense that for sampling from a very large collection with  $\frac{m}{N} = p$ , sampling with or without replacement are practically the same.

### 2.1.6 Poisson Distribution $\text{Poi}(\lambda)$

The Poisson Distribution is an approximation of the Binomial Distribution  $\text{Bin}(n, p)$ , where  $n$  is large,  $p$  is small, and  $np$  is of moderate magnitude. So let  $n \rightarrow \infty$ ,  $p \rightarrow 0$  in such a manner that  $np \rightarrow \lambda \in (0, \infty)$ , and approximate the pmf of  $\text{Bin}(n, p)$  for such  $n, p$  as follows:

$$\begin{aligned}
f_X(x; n, p) &= \frac{n(n-1) \cdots (n-x+1)}{x!} p^x (1-p)^{n-x} \\
&= \frac{(np)^x}{x!} \left[ \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x-1}{n}\right) \right] [1-p]^{-x} \left[1 - \frac{np}{n}\right]^n \\
&:= A_{1n} A_{2n} A_{3n} A_{4n},
\end{aligned}$$

where  $A_{1n} \rightarrow \lambda^x/x!$ ,  $A_{2n} \rightarrow 1$ ,  $A_{3n} \rightarrow 1$  and since  $\left(1 - \frac{a}{n}\right)^n \rightarrow e^{-a}$ , which extends to  $\left(1 - \frac{a_n}{n}\right)^n \rightarrow e^{-a}$ , if  $a_n \rightarrow a$ , we see that  $A_{4n} \rightarrow e^{-\lambda}$ . Thus the pmf of  $\text{Bin}(n, p)$ , as  $n \rightarrow \infty$

and  $p \rightarrow 0$  in the above manner, tends to the following:

$$f_X(x; n, p) \rightarrow f_X(x) = \frac{\lambda^x}{x!} \cdot 1 \cdot 1 \cdot e^{-\lambda} = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

This limiting pmf of  $\text{Bin}(n, p)$  as  $n \rightarrow \infty, p \rightarrow 0$  so that  $np \rightarrow \lambda \in (0, \infty)$  is the pmf of the Poisson distribution  $\text{Poi}(\lambda)$ :

$$f_X(x) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

and

$$\sum_{x=0}^{\infty} f_X(x) = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1.$$

The mean and the variance of  $\text{Poi}(\lambda)$  are

$$\begin{aligned} E[X] &= \sum_{x=0}^{\infty} x e^{-\lambda} \frac{\lambda^x}{x!} = \sum_{x=1}^{\infty} e^{-\lambda} \frac{\lambda^x}{(x-1)!} \\ &= \lambda \sum_{y=0}^{\infty} e^{-\lambda} \frac{\lambda^y}{y!} = \lambda \cdot 1 = \lambda \text{ and,} \\ E[X(X-1)] &= \sum_{x=0}^{\infty} x(x-1) e^{-\lambda} \frac{\lambda^x}{x!} = \sum_{x=2}^{\infty} e^{-\lambda} \frac{\lambda^x}{(x-2)!} \\ &= \lambda^2 \sum_{y=0}^{\infty} e^{-\lambda} \frac{\lambda^y}{y!} = \lambda^2 \cdot 1 = \lambda^2, \text{ so} \\ \text{Var}[X] &= E[X^2] - (E[X])^2 = E[X(X-1)] + E[X] - (E[X])^2 \\ &= \lambda^2 + \lambda - \lambda^2 = \lambda. \end{aligned}$$

*Some Properties of  $\text{Poi}(\lambda)$ .*

- Let  $X_1, X_2$  be independent,  $X_i \sim \text{Poi}(\lambda_i)$ . Then  $X = X_1 + X_2 \sim \text{Poi}(\lambda_1 + \lambda_2)$ , by Eq. (12b) of [Chapter 1](#),

$$\begin{aligned} f_X(x) &= \sum_{x_1=0}^x f_{X_1}(x_1) f_{X_2}(x-x_1) = e^{-(\lambda_1+\lambda_2)} \sum_{x_1=0}^x \frac{\lambda_1^{x_1}}{x_1!} \frac{\lambda_2^{x-x_1}}{(x-x_1)!} \\ &= e^{-(\lambda_1+\lambda_2)} \cdot \frac{(\lambda_1+\lambda_2)^x}{x!} \sum_{x_1=0}^x \binom{x}{x_1} \left(\frac{\lambda_1}{\lambda_1+\lambda_2}\right)^{x_1} \left(\frac{\lambda_2}{\lambda_1+\lambda_2}\right)^{x-x_1} \\ &= e^{-(\lambda_1+\lambda_2)} \cdot \frac{(\lambda_1+\lambda_2)^x}{x!} \cdot 1, \end{aligned}$$

which is the pmf of  $\text{Poi}(\lambda_1 + \lambda_2)$ . More generally,  $X_1, \dots, X_k$  independent with  $X_i \sim \text{Poi}(\lambda_i)$  implies  $\sum_{i=1}^k X_i \sim \text{Poi}\left(\sum_{i=1}^k \lambda_i\right)$ .

2. If  $X_1, X_2$  are independent with  $X_i \sim Poi(\lambda_i)$ , then conditionally, given

$$X_1 + X_2 = n, X_1 \sim Bin(n, \lambda_1/(\lambda_1 + \lambda_2)).$$

*Proof.*

$$\begin{aligned} P[X_1 = x | X_1 + X_2 = n] &= \frac{P[X_1 = x, X_2 = n - x]}{P[X_1 + X_2 = n]} \\ &= \frac{e^{-(\lambda_1 + \lambda_2)} \frac{\lambda_1^x}{x!} \frac{\lambda_2^{n-x}}{(n-x)!}}{e^{-(\lambda_1 + \lambda_2)} \frac{(\lambda_1 + \lambda_2)^n}{n!}} \\ &= \frac{n!}{x!(n-x)!} \left( \frac{\lambda_1}{\lambda_1 + \lambda_2} \right)^x \left( \frac{\lambda_2}{\lambda_1 + \lambda_2} \right)^{n-x}, \end{aligned}$$

which is the pmf of  $Bin(n, \lambda_1/(\lambda_1 + \lambda_2))$ .  $\square$

More generally, if  $X_i \sim Poi(\lambda_i)$   $i = 1, \dots, k$  are independent, then conditionally, given  $X_1 + \dots + X_k = n$ ,  $(X_1, \dots, X_k) \sim Mult(n; p_1, \dots, p_k)$  with  $p_i = \lambda_i/(\lambda_1 + \dots + \lambda_k)$ .

3. If  $N \sim Poi(\lambda)$  and conditionally, given  $N = n$ ,  $X_1, \dots, X_n$  are independent  $Bernoulli(p)$ , then  $X_1 + \dots + X_N = X \sim Poi(\lambda p)$ .

*Proof.*

$$\begin{aligned} P[X = x] &= \sum_{n=x}^{\infty} P[N = n, X_1 + \dots + X_n = x] \\ &= \sum_{n=x}^{\infty} \frac{e^{-\lambda} \lambda^n}{n!} \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \\ &= e^{-\lambda} \cdot \frac{(\lambda p)^x}{x!} \sum_{n=x}^{\infty} \frac{\{\lambda(1-p)\}^{n-x}}{(n-x)!} \\ &= e^{-\lambda} \frac{(\lambda p)^x}{x!} \cdot e^{\lambda(1-p)} = e^{-\lambda p} \frac{(\lambda p)^x}{x!}. \end{aligned}$$

$\square$

**Example 2.1.3.** An insurance agent has sold fire insurance policies to 200 homeowners in a town. If the probability of a fire in a house during a year is 1/250 in this town, what is the probability that this agent will have to handle two or more claims in a year?

*Solution.* The number of policies  $n = 200$  (large) and the probability of a fire is  $p = 1/250$  (small), so  $\lambda = np = 0.8$ . Assuming independence,  $X = \text{number of claims} = \text{number of fires}$ , is approximately  $Poi(\lambda = 0.8)$ . Hence

$$\begin{aligned} P[X \geq 2] &= 1 - (P[X = 0] + P[X = 1]) = 1 - e^{-0.8}[1 + 0.8] \\ &= 1 - 0.4493 \cdot 1.8 = 1 - 0.8088 = 0.1912. \end{aligned}$$

**Example 2.1.4.** If the number of tropical storms on the Gulf of Mexico during September follows a Poisson distribution with mean  $\lambda = 5$  and if each of these storms can become a hurricane with probability  $p = 1/4$ , what is the probability of at least one hurricane in September?

*Solution.*  $N = \text{Number of storms} \sim Poi(\lambda = 5)$ , and  $X_1, \dots, X_N$  given  $N$  are iid  $Bernoulli(p = 1/4)$ , where  $X_i = 1$  or 0 according as the  $i$ th storm becomes a hurricane or not. Hence  $X = \text{number of hurricanes} = X_1 + \dots + X_N \sim Poi(\lambda p = 1.25)$ , and so

$$\begin{aligned}
P[\text{At least one hurricane}] &= P[X \geq 1] \\
&= 1 - P[X = 0] = 1 - e^{-1.25} \\
&= 1 - 0.2865 = 0.7135.
\end{aligned}$$

## 2.2 Continuous Distributions

### 2.2.1 The Gamma and Beta Functions

These two functions are essential in connection with many distributions and their properties discussed here.

**Definition 2.2.1.** The gamma function is defined as

$$\Gamma(\alpha) = \int_0^\infty e^{-u} u^{\alpha-1} du, \quad \alpha > 0$$

and the beta function is defined as

$$Be(\alpha_1, \alpha_2) = \int_0^1 t^{\alpha_1-1} (1-t)^{\alpha_2-1} dt, \quad \alpha_1, \alpha_2 > 0.$$

*Properties of the Gamma and Beta Functions.*

(i)  $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$  for  $\alpha > 0$ , so for a positive integer  $\alpha$ ,  $\Gamma(\alpha + 1) = \alpha!$ .

*Proof.* Integrating by parts,

$$\begin{aligned}
\Gamma(\alpha + 1) &= \int_0^\infty e^{-u} u^\alpha du = - \int_0^\infty u^\alpha d(e^{-u}) \\
&= \int_0^\infty e^{-u} d(u^\alpha) = \alpha \int_0^\infty e^{-u} u^{\alpha-1} du \\
&= \alpha \Gamma(\alpha),
\end{aligned}$$

because  $e^{-u} u^\alpha = 0$  at  $u = 0$  and  $\rightarrow 0$  as  $u \rightarrow \infty$ . □

(ii)  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ .

*Proof.*

$$\begin{aligned}
\left\{ \Gamma\left(\frac{1}{2}\right) \right\}^2 &= \int_0^\infty \int_0^\infty e^{-u-v} u^{-1/2} v^{-1/2} du dv \\
&= \int_0^\infty \int_0^1 e^{-s} (st)^{-1/2} \{s(1-t)\}^{-1/2} \cdot s ds dt,
\end{aligned}$$

with  $u + v = s$ ,  $u/(u + v) = t$ , so that  $du dv = s ds dt$ . Next let  $\sin^2 \theta = t$ , so that  $dt = 2 \sin \theta \cos \theta d\theta$  and  $\{t(1-t)\}^{1/2} = \sin \theta \cos \theta$ . Hence

$$\begin{aligned}
\left\{ \Gamma\left(\frac{1}{2}\right) \right\}^2 &= \int_0^\infty e^{-s} ds \cdot \int_0^1 \{t(1-t)\}^{-1/2} dt \\
&= 1 \cdot \int_0^{\pi/2} \frac{2 \sin \theta \cos \theta}{\sin \theta \cos \theta} d\theta \\
&= 2 \cdot \frac{\pi}{2} = \pi. \quad \square
\end{aligned}$$

(iii)  $Be(\alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(\alpha_1+\alpha_2)}$ . This will follow as part of the proof of [Proposition 2.2.1](#).

We now introduce some important continuous distributions and examine some of their properties.

### 2.2.2 Uniform Distribution $Unif(a, b)$

**Definition 2.2.2.**  $X$  is a  $Unif(a, b)$  rv if  $f_X(x) = I_{[a,b]}(x)/(b-a)$ .

Of course,  $\int_a^b (b-a)^{-1} dx = (b-a)^{-1}(b-a) = 1$ .

$$\begin{aligned}\mathbb{E}[X] &= \int_a^b \frac{x}{b-a} dx = \frac{1}{b-a} \cdot \frac{b^2 - a^2}{2} = \frac{a+b}{2}, \\ \mathbb{E}[X^2] &= \int_a^b \frac{x^2}{b-a} dx = \frac{b^3 - a^3}{3(b-a)} = \frac{a^2 + ab + b^2}{3}, \text{ and} \\ \text{Var}[X] &= \frac{a^2 + ab + b^2}{3} - \frac{a^2 + 2ab + b^2}{4} = \frac{(b-a)^2}{12}.\end{aligned}$$

In particular, if  $X \sim Unif(0, 1)$ , then  $f_X(x) = I_{[0,1]}(x)$ ,  $\mathbb{E}[X] = \frac{1}{2}$ , and  $\text{Var}[X] = \frac{1}{12}$ .

### 2.2.3 Gamma and Beta distributions

**Definition 2.2.3.**  $X$  is a  $Gamma(\alpha, \beta)$  rv if the pdf of  $X$  is

$$f_X(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} e^{-x/\beta} x^{\alpha-1} I_{(0,\infty)}(x), \quad \alpha > 0, \quad \beta > 0.$$

The mean and variance of a  $Gamma(\alpha, \beta)$  rv are

$$\begin{aligned}\mathbb{E}[X] &= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^\infty x e^{-x/\beta} x^{\alpha-1} dx = \frac{\beta}{\Gamma(\alpha)} \int_0^\infty e^{-u} u^\alpha du \\ &= \frac{\beta}{\Gamma(\alpha)} \Gamma(\alpha+1) = \alpha\beta, \\ \text{Var}[X] &= \mathbb{E}[X^2] - (\alpha\beta)^2 = \frac{\beta^2}{\Gamma(\alpha)} \int_0^\infty e^{-u} u^{\alpha+1} du - (\alpha\beta)^2 \\ &= \frac{\beta^2}{\Gamma(\alpha)} \Gamma(\alpha+2) - (\alpha\beta)^2 \\ &= \alpha(\alpha+1)\beta^2 - \alpha^2\beta^2 = \alpha\beta^2.\end{aligned}$$

**Proposition 2.2.1.** If  $X_i$ ,  $i = 1, 2$  are independent  $Gamma(\alpha_i, \beta)$  rv's, then  $X_1 + X_2$  is a  $Gamma(\alpha_1 + \alpha_2, \beta)$  rv.

*Proof.* By Eq. (12a) in Chapter 1,

$$\begin{aligned}
 f_{X_1+X_2}(z) &= \frac{1}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} \int_0^z e^{-x/\beta} x^{\alpha_1-1} e^{-(z-x)/\beta} (z-x)^{\alpha_2-1} dx \\
 &= \frac{e^{-z/\beta} z^{\alpha_1+\alpha_2-1}}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} \int_0^z \left(\frac{x}{z}\right)^{\alpha_1-1} \left(1 - \frac{x}{z}\right)^{\alpha_2-1} \frac{dx}{z} \\
 &= \frac{e^{-z/\beta} z^{\alpha_1+\alpha_2-1}}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} \int_0^1 t^{\alpha_1-1} (1-t)^{\alpha_2-1} dt \\
 &= \frac{e^{-z/\beta} z^{\alpha_1+\alpha_2-1}}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} \cdot Be(\alpha_1, \alpha_2), \quad z > 0.
 \end{aligned}$$

But we must have

$$\begin{aligned}
 1 &= \int_0^\infty f_{X_1+X_2}(z) dz \\
 &= \frac{Be(\alpha_1, \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} \int_0^\infty e^{-z/\beta} z^{\alpha_1+\alpha_2-1} dz \\
 &= \frac{Be(\alpha_1, \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \Gamma(\alpha_1 + \alpha_2).
 \end{aligned}$$

Hence  $Be(\alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(\alpha_1+\alpha_2)}$  as stated earlier, and therefore,

$$f_{X_1+X_2}(z) = \frac{1}{\Gamma(\alpha_1 + \alpha_2)\beta^{\alpha_1+\alpha_2}} e^{-z/\beta} z^{\alpha_1+\alpha_2-1} I_{0,\infty}(z).$$

□

By induction, if  $X_i$ ,  $i = 1, \dots, k$  are independent  $Gamma(\alpha_i, \beta)$ , then  $\sum_{i=1}^k X_i \sim Gamma\left(\sum_{i=1}^k \alpha_i, \beta\right)$ .

**Definition 2.2.4.**  $T$  is a  $Be(\alpha_1, \alpha_2)$  rv if the pdf of  $T$  is

$$f_T(t) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} t^{\alpha_1-1} (1-t)^{\alpha_2-1} I_{(0,1)}(t), \quad \alpha_1, \alpha_2 > 0.$$

The mean and variance of a  $Beta(\alpha_1, \alpha_2)$  rv are

$$\begin{aligned}
 E[T] &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \int_0^1 t^{\alpha_1} (1-t)^{\alpha_2-1} dt \\
 &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} Be(\alpha_1 + 1, \alpha_2) \\
 &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \cdot \frac{\Gamma(\alpha_1 + 1)\Gamma(\alpha_2)}{\Gamma(\alpha_1 + \alpha_2 + 1)} \\
 &= \frac{\alpha_1}{\alpha_1 + \alpha_2},
 \end{aligned}$$

$$\text{Var}[T] = E[T^2] - \left(\frac{\alpha_1}{\alpha_1 + \alpha_2}\right)^2$$

$$\begin{aligned}
&= \frac{\alpha_1(\alpha_1 + 1)}{(\alpha_1 + \alpha_2)(\alpha_1 + \alpha_2 + 1)} - \left( \frac{\alpha_1}{\alpha_1 + \alpha_2} \right)^2 \\
&= \frac{\alpha_1 \alpha_2}{(\alpha_1 + \alpha_2)^2 (\alpha_1 + \alpha_2 + 1)}.
\end{aligned}$$

**Proposition 2.2.2.** If  $X_i \sim \text{Gamma}(\alpha_i, \beta)$ ,  $i = 1, 2$  are independent, then  $Z = X_1 + X_2$  and  $T = X_1/(X_1 + X_2)$  are independent,  $Z \sim \text{Gamma}(\alpha_1 + \alpha_2, \beta)$  and  $T \sim \text{Be}(\alpha_1, \alpha_2)$ .

*Proof.* Rewrite the transformation as:  $X_1 = TZ$ ,  $X_2 = (1 - T)Z$ . Then the Jacobian of the transformation is  $|J| = z$ , and

$$\begin{aligned}
f_{Z,T}(z, t) &= \frac{1}{\Gamma(\alpha_1)\Gamma(\alpha_2)\beta^{\alpha_1+\alpha_2}} e^{-z/\beta} (tz)^{\alpha_1-1} \{(1-t)z\}^{\alpha_2-1} z I_{(0,\infty)}(z) I_{(0,1)}(t) \\
&= \frac{e^{-z/\beta} z^{\alpha_1+\alpha_2-1}}{\Gamma(\alpha_1+\alpha_2)\beta^{\alpha_1+\alpha_2}} I_{(0,\infty)}(z) \cdot \frac{\Gamma(\alpha_1+\alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} t^{\alpha_1-1} (1-t)^{\alpha_2-1} I_{(0,1)}(t),
\end{aligned}$$

as was to be shown.  $\square$

## 2.2.4 Exponential Distribution

**Definition 2.2.5.** A rv  $X$  is distributed as  $\text{Exp}(\theta)$  if  $X$  has pdf

$$f_X(x) = \frac{1}{\theta} e^{-x/\theta}, \quad \theta > 0,$$

and cdf

$$F_X(t) = \frac{1}{\theta} \int_0^t e^{-x/\theta} dx = 1 - e^{-t/\theta}.$$

Clearly,  $\text{Exp}(\theta)$  is a special case of  $\text{Gamma}(\alpha, \beta)$  with  $\alpha = 1$  and  $\beta = \theta$ . Hence,  $E[X] = \theta$  and  $\text{Var}[X] = \theta^2$ .

The exponential distribution shares with geometric distribution, the memoryless property:

$$P[X > t+s | X > t] = \frac{P[X > t+s]}{P[X > t]} = \frac{e^{-(t+s)/\theta}}{e^{-t/\theta}} = e^{-s/\theta} = P[X > s].$$

## 2.2.5 Normal Distribution $N(\mu, \sigma^2)$

**Definition 2.2.6.**  $X$  is an  $N(\mu, \sigma^2)$  rv if  $X$  has pdf

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}, \quad -\infty < x < \infty.$$

We first verify that  $f_X(x)$  is indeed a pdf. Since  $f_X(x) \geq 0$ , we only need to check that  $\int_{-\infty}^{\infty} f_X(x) dx = 1$ . By symmetry,

$$\begin{aligned}
\int_{-\infty}^{\infty} f_X(x) dx &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-z^2/2} dz = \frac{2}{\sqrt{2\pi}} \int_0^{\infty} e^{-z^2/2} dz \\
&= \frac{2}{\sqrt{2\pi}} \int_0^{\infty} e^{-u} \cdot \frac{1}{\sqrt{2}} u^{-1/2} du \quad \left( \text{with } u = \frac{z^2}{2}, \quad dz = (2u)^{-1/2} du \right)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{\pi}} \int_0^\infty e^{-u} u^{-1/2} du = \frac{1}{\sqrt{\pi}} \cdot \Gamma\left(\frac{1}{2}\right) \\
&= 1, \text{ since } \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}.
\end{aligned}$$

The mean and variance of  $N(\mu, \sigma^2)$  are

$$\begin{aligned}
E[X] &= \int_{-\infty}^\infty x f_X(x) dx = \int_{-\infty}^\infty (x - \mu) f_X(x) dx + \mu \cdot 1 \\
&= \int_{-\infty}^\infty \sigma z \cdot \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz + \mu \\
&= \mu, \text{ because } \int_{-\infty}^\infty z e^{-z^2/2} dz = 0 \text{ by symmetry.} \\
\text{Var}[X] &= E[(X - \mu)^2] \\
&= \int_{-\infty}^\infty (x - \mu)^2 \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{2\sigma^2}{\sqrt{2\pi}} \int_0^\infty z^2 e^{-z^2/2} dz \\
&= \frac{2\sigma^2}{\sqrt{2\pi}} \int_0^\infty 2ue^{-u} \cdot (2u)^{-1/2} du \quad \left( \text{with } u = \frac{z^2}{2} \text{ as above} \right) \\
&= \frac{2\sigma^2}{\sqrt{\pi}} \int_0^\infty e^{-u} u^{1/2} du = \frac{2\sigma^2}{\sqrt{\pi}} \cdot \Gamma\left(\frac{3}{2}\right) \\
&= \frac{2\sigma^2}{\sqrt{\pi}} \cdot \frac{1}{2} \Gamma\left(\frac{1}{2}\right) = \sigma^2.
\end{aligned}$$

For this reason,  $N(\mu, \sigma^2)$  is called a normal distribution with mean  $\mu$  and variance  $\sigma^2$ . For  $\mu = 0$  and  $\sigma = 1$ ,  $N(0,1)$  is called the standard normal distribution.

**Example 2.2.1.** The life of a tire on freeway is a Normal rv with mean  $\mu = 25,000$  miles and sd  $\sigma = 1,000$  miles. If a car goes on a coast-to-coast round-trip in the United States, which we can assume is 6000 miles, what is the probability that none of the four tires will need replacement during the 6000 miles trip, given that the tires were already used for 17,500 miles?

*Solution.* Note that

$$\begin{aligned}
P[\min(X_1, \dots, X_4) > 23,500] &= \{P[X_1 > 23,500]\}^4 \\
&= \left\{ P\left[\frac{X_1 - \mu}{\sigma} = Z > \frac{23,500 - 25,000}{1000}\right] \right\}^4 \\
&= \{P[Z > -1.5]\}^4 = \{P[Z < 1.5]\}^4 \\
&= \{\Phi(1.5)\}^4 = \{1 - 0.0668\}^4 = 0.7584,
\end{aligned}$$

using the normal table.

Strictly speaking, in this problem, we should have considered

$$P[\min(X_1, \dots, X_4) > 23,500 | X_i > 17,500, i = 1, \dots, 4],$$

but

$$P[X_i > 17,500] = P\left[Z > \frac{17,500 - 25,000}{1000}\right] = P[Z > -7.5] \approx 1,$$

so this conditionality does not matter.

*pdf and cdf of  $N(0, 1)$ .* Let  $Z \sim N(0, 1)$ . Then the pdf of  $Z$  is

$$f_Z(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

which we shall denote by  $\varphi(x)$ . The cdf of  $Z$  is  $F_Z(x) = \int_{-\infty}^x \varphi(t) dt$ , which we shall denote by  $\Phi(x)$ . Since  $\varphi(x)$  is symmetric about 0, we have

$$\Phi(-x) = 1 - \Phi(x) \text{ and } P[|Z| > x] = 2[1 - \Phi(x)].$$

The cdf  $\Phi(x)$  cannot be calculated analytically. There are tables for  $\Phi(x)$ . However, the following are useful upper and lower bounds for  $1 - \Phi(x)$  for all  $x > 0$ ,

$$x^{-1}(1 - x^{-2})\varphi(x) < 1 - \Phi(x) < x^{-1}\varphi(x).$$

By a change of variable  $y = x^2/2$ ,

$$1 - \Phi(x) = \int_x^\infty \varphi(t) dt = (1/2)\pi^{-1/2} \int_{x^2/2}^\infty y^{-1/2} e^{-y} dy.$$

For notational convenience use  $c = x^2/2$ . The upper bound is obtained by noting that

$$(1/2)\pi^{-1/2} \int_c^\infty y^{-1/2} e^{-y} dy < (1/2)\pi^{-1/2} c^{-1/2} \int_c^\infty e^{-y} dy = (1/2)\pi^{-1/2} c^{-1/2} e^{-c} = x^{-1}\varphi(x).$$

To get the lower bound, use the integration by parts formula

$$\begin{aligned} (1/2)\pi^{-1/2} \int_c^\infty y^{-1/2} e^{-y} dy &= (1/2)\pi^{-1/2} \left[ c^{-1/2} e^{-c} - \frac{1}{2} \int_c^\infty y^{-3/2} e^{-y} dy \right] \\ &> (1/2)\pi^{-1/2} \left[ c^{-1/2} e^{-c} - \frac{1}{2} c^{-3/2} \int_c^\infty e^{-c} \right] \\ &= (1/2)\pi^{-1/2} \left[ c^{-1/2} e^{-c} - \frac{1}{2} c^{-3/2} e^{-c} \right] \\ &= (1/2)\pi^{-1/2} c^{-1/2} \left[ 1 - \frac{1}{2} c^{-1} \right] e^{-c} \\ &= x^{-1}(1 - x^{-2})\varphi(x). \end{aligned}$$

*Normal Approximation to the Binomial Distribution: De Moivre-Laplace Theorem.* If  $S_n = X_1 + \dots + X_n$ , where  $X_1, \dots, X_n$  are independent  $Bernoulli(p)$ , that is,  $S_n \sim Bin(n, p)$ , then for large  $n$  and for any  $a < b$ ,

$$P\left[a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b\right] \approx \Phi(b) - \Phi(a).$$

This is a special case of the Central Limit Theorem which holds for sums of iid rv's  $S_n = X_1 + \dots + X_n$  with finite variance.

*Some Properties of  $N(\mu, \sigma^2)$ .*

1. If  $X \sim N(\mu, \sigma^2)$ , then  $a + bX \sim N(a + b\mu, b^2\sigma^2)$  follows by using Eq. (10) in Chapter 1. In particular,  $Z = (X - \mu)/\sigma \sim N(0, 1)$ .

2. If  $X_i \sim N(\mu_i, \sigma_i^2)$ ,  $i = 1, 2$  are independent, then  $X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$ .

*Proof.*  $Z_i = (X_i - \mu_i)/\sigma_i$ ,  $i = 1, 2$  are independent  $N(0, 1)$ , and

$$X_1 + X_2 = (\mu_1 + \mu_2) + (\sigma_1 Z_1 + \sigma_2 Z_2). \text{ Transform } \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} \sigma_1 & \sigma_2 \\ \sigma_2 & -\sigma_1 \end{pmatrix} \cdot \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix}. \text{ Then}$$

$\left| \det \begin{pmatrix} \sigma_1 & \sigma_2 \\ \sigma_2 & -\sigma_1 \end{pmatrix} \right| = \sigma_1^2 + \sigma_2^2$ , and  $\begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} = \frac{1}{\sqrt{\sigma_1^2 + \sigma_2^2}} \begin{pmatrix} \sigma_1 & \sigma_2 \\ \sigma_2 & -\sigma_1 \end{pmatrix} \cdot \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$ , and the joint pdf of  $Y_1 = \sigma_1 Z_1 + \sigma_2 Z_2$  and  $Y_2 = \sigma_2 Z_1 - \sigma_1 Z_2$  is

$$\begin{aligned} f_{Y_1, Y_2}(y_1, y_2) &= \frac{1}{2\pi(\sigma_1^2 + \sigma_2^2)} \exp \left[ -\frac{1}{2(\sigma_1^2 + \sigma_2^2)^2} \left\{ (\sigma_1 y_1 + \sigma_2 y_2)^2 + (\sigma_2 y_1 - \sigma_1 y_2)^2 \right\} \right] \\ &= \frac{1}{2\pi(\sigma_1^2 + \sigma_2^2)} \exp \left[ -\frac{1}{2(\sigma_1^2 + \sigma_2^2)^2} \left\{ (\sigma_1^2 + \sigma_2^2)y_1^2 + (\sigma_1^2 + \sigma_2^2)y_2^2 \right\} \right] \\ &= \frac{1}{2\pi(\sigma_1^2 + \sigma_2^2)} \exp \left[ -\frac{y_1^2}{2(\sigma_1^2 + \sigma_2^2)} - \frac{y_2^2}{2(\sigma_1^2 + \sigma_2^2)} \right] \end{aligned}$$

Thus  $Y_1, Y_2$  are independent  $N(0, \sigma_1^2 + \sigma_2^2)$ , so

$$\begin{aligned} X_1 + X_2 &= (\mu_1 + \mu_2) + (\sigma_1 Z_1 + \sigma_2 Z_2) \\ &= (\mu_1 + \mu_2) + Y_1 \\ &\sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2). \end{aligned}$$

□

By induction, if  $X_i \sim N(\mu_i, \sigma_i^2)$ ,  $i = 1, \dots, k$  are mutually independent, then

$\sum_{i=1}^k X_i \sim N\left(\sum_{i=1}^k \mu_i, \sum_{i=1}^k \sigma_i^2\right)$ . In particular, if  $X_1, \dots, X_n$  are independent  $N(\mu, \sigma^2)$ , then  $\sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2)$ , and  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \sigma^2/n)$ .

3. If  $Z \sim N(0, 1)$ , then  $W = Z^2$  has pdf

$$f_W(w) = \frac{1}{2^{1/2} \Gamma\left(\frac{1}{2}\right)} e^{-w/2} w^{-1/2} I_{(0, \infty)}(w). \quad (3)$$

*Proof.*  $W = Z^2$  has cdf (by Eq. (11b) in Chapter 1)

$$F_{Z^2}(w) = P[-\sqrt{w} \leq Z \leq \sqrt{w}] = 2 \cdot P[0 \leq Z \leq \sqrt{w}]$$

$$= 2 \int_0^{\sqrt{w}} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz,$$

so

$$f_{Z^2}(w) = \sqrt{\frac{2}{\pi}} e^{-w/2} \cdot \frac{1}{2} w^{-1/2} = \frac{1}{2^{1/2} \Gamma\left(\frac{1}{2}\right)} e^{-w/2} w^{-1/2}, \quad w > 0.$$

□

## 2.2.6 Chi-Square Distribution

If  $Z \sim N(0, 1)$ , then from the pdf of  $Z^2$  derived in Eq. (3) we see that  $Z^2$  is a  $\text{Gamma}\left(\frac{1}{2}, 2\right)$  rv and from the additive property of independent Gamma rv's with common  $\beta$  in [Proposition 2.2.1](#), the following proposition is an immediate consequence.

**Proposition 2.2.3.** *If  $Z_1, \dots, Z_k$  are independent  $N(0, 1)$ , then  $W = \sum_{i=1}^k Z_i^2$  is a  $\text{Gamma}(k/2, 2)$  rv with pdf:*

$$f_W(w) = \frac{1}{2^{k/2} \Gamma\left(\frac{k}{2}\right)} e^{-w/2} w^{k/2-1} I_{(0, \infty)}(w).$$

This rv is known as  $\chi^2$  with  $k$  degrees of freedom (df), denoted by  $\chi_k^2$ . Clearly, the sum of independent  $\chi^2$  rv's is a  $\chi^2$  rv, the df of which is the sum of the df's of its components.

**Proposition 2.2.4.** *If  $\chi_{k_1}^2, \dots, \chi_{k_n}^2$  are mutually independent, then by the additive property of independent Gamma rv's derived in [Proposition 2.2.1](#),  $\chi_{k_1}^2 + \dots + \chi_{k_n}^2$  is distributed as  $\chi_{k_1+\dots+k_n}^2$ .*

Now suppose that  $X \sim N(\mu, 1)$ , then by Eq. (11a) of [Chapter 1](#),

$$\begin{aligned} f_{X^2}(w) &= \frac{1}{2\sqrt{w}} \cdot \frac{1}{\sqrt{2\pi}} \left[ e^{-\frac{1}{2}(\sqrt{w}-\mu)^2} + e^{-\frac{1}{2}(-\sqrt{w}-\mu)^2} \right] \\ &= \frac{1}{2\sqrt{w}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\mu^2 - \frac{1}{2}w} \left[ e^{\mu\sqrt{w}} + e^{-\mu\sqrt{w}} \right]. \end{aligned}$$

Using the properties of Gamma function and remembering  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ ,

$$\begin{aligned} e^a + e^{-a} &= 2 \sum_{r=0}^{\infty} \frac{a^{2r}}{(2r)!} = 2 \sum_{r=0}^{\infty} \frac{(a^2)^r}{2^r r! [1 \cdot 3 \cdots (2r-1)]} \\ &= \sqrt{\pi} \sum_{r=0}^{\infty} \frac{(a^2)^r / 2^{2r-1}}{r! \left[ \frac{1}{2} \cdot \frac{3}{2} \cdots \left(r - \frac{3}{2}\right) \cdot \left(r - \frac{1}{2}\right) \right] \Gamma\left(\frac{1}{2}\right)} \\ &= \sqrt{\pi} \sum_{r=0}^{\infty} \frac{(a^2)^r / 2^{2r-1}}{r! \Gamma\left(r + \frac{1}{2}\right)}. \end{aligned}$$

We thus arrive at the *Noncentral  $\chi^2$  distribution with 1 df* and noncentrality parameter  $\frac{1}{2}\mu^2$ :

$$\begin{aligned} f_{X^2}(w) &= \frac{1}{2\sqrt{w}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\mu^2 - \frac{1}{2}w} \cdot \sqrt{\pi} \sum_{r=0}^{\infty} \frac{\left(\frac{1}{2}\mu^2\right)^r w^{r/2} r^{2r-1}}{r! \Gamma\left(r + \frac{1}{2}\right)} \\ &= \sum_{r=0}^{\infty} \left( \frac{e^{-\frac{1}{2}\mu^2} \left(\frac{1}{2}\mu^2\right)^r}{r!} \right) \cdot \frac{1}{\Gamma\left(\frac{2r+1}{2}\right) 2^{\frac{2r+1}{2}}} e^{-w/2} w^{\frac{2r+1}{2}-1}. \end{aligned}$$

Let  $N$  be a Poisson rv with mean  $\frac{1}{2}\mu^2$ . Then the  $r$ th term in the above sum is

$$f_N\left(r, \frac{1}{2}\mu^2\right) \cdot f_{\chi_{1+2r}^2}(w).$$

Hence,  $f_{X^2}(w) = E[f_{\chi_{1+2N}^2}(w)]$ , where  $N$  is a  $Poi\left(\frac{1}{2}\mu^2\right)$  rv. This is the pdf of a noncentral  $\chi^2$  rv with 1 df and noncentrality parameter  $\frac{1}{2}\mu^2$  denoted by  $\chi_1^2\left(\frac{1}{2}\mu^2\right)$ .

### 2.2.7 Sampling From a Normal Distribution, Sample Mean and Sample Variance

Let  $X_1, \dots, X_n$  be  $n$  independent rv's, each distributed as  $N(\mu, \sigma^2)$ . Such a collection of rv's will be called "a random sample of  $n$  observations from the normal population  $N(\mu, \sigma^2)$ ." The concepts of population and sample will be discussed in [Chapter 4](#).

We have seen that  $\mu$  is the mean and  $\sigma^2$  is the variance of  $N(\mu, \sigma^2)$ . The analogs of  $\mu$  and  $\sigma^2$  in the sample  $(X_1, \dots, X_n)$  are, respectively, the sample mean and sample variance

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \text{ and } s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

(The reason for the factor  $\frac{1}{n-1}$  instead of  $\frac{1}{n}$  in  $s^2$  will be made clear later in [Chapter 5](#)). In the following proposition we derive the joint distribution of  $\bar{X}$  and  $s^2$ .

**Proposition 2.2.5.** *Let  $\bar{X}$  and  $s^2$  be the sample mean and sample variance, respectively, in a random sample  $(X_1, \dots, X_n)$  from  $N(\mu, \sigma^2)$ . Then*

- (i)  $\bar{X} \sim N(\mu, \sigma^2/n)$ ,
- (ii)  $(n-1)s^2/\sigma^2 \sim \chi_{n-1}^2$ , and
- (iii)  $\bar{X}$  and  $s^2$  are independent.

*Proof.* Let  $W_i = X_i - \mu$ . Then  $W_1, \dots, W_n$  are iid  $N(0, \sigma^2)$ , and letting  $\mathbf{W}^\top = (W_1, \dots, W_n)$ ,  $\mathbf{w}^\top = (w_1, \dots, w_n)$ , the joint pdf of  $(W_1, \dots, W_n)$  is expressed as

$$f_{\mathbf{W}}(\mathbf{w}) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left[-\frac{1}{2\sigma^2} \mathbf{w}^\top \mathbf{w}\right]$$

Transform  $\mathbf{Y} = \mathbf{A}\mathbf{W}$ , where  $\mathbf{A} = ((a_{ij}))$  with  $\mathbf{a}_i^\top = (a_{i1}, \dots, a_{in})$ ,  $i = 1, \dots, n$  satisfying

$$\mathbf{a}_1^\top = n^{-1/2}(1, \dots, 1), \mathbf{a}_i^\top \mathbf{a}_i = 1, \text{ and } \mathbf{a}_i^\top \mathbf{a}_j = 0 \text{ for all } i \neq j$$

so that  $\mathbf{A}$  is orthonormal. (In particular,

$$\mathbf{a}_j^\top = \{(j-1)j\}^{-1/2}(1, \dots, 1, j-1, 0, \dots, 0), \quad j = 2, \dots, n$$

will do.) Then

$$Y_1 = \mathbf{a}_1^\top \mathbf{W} = n^{-1/2} \sum_{i=1}^n W_i = \sqrt{n}\bar{W} = \sqrt{n}(\bar{X} - \mu),$$

where  $\bar{W} = n^{-1} \sum_{i=1}^n W_i$ , and  $\mathbf{Y}^\top \mathbf{Y} = \mathbf{W}^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{W} = \mathbf{W}^\top \mathbf{W}$ . Hence

$$\sum_{i=2}^n Y_i^2 = \mathbf{Y}^\top \mathbf{Y} - Y_1^2 = \mathbf{W}^\top \mathbf{W} - (\sqrt{n}\bar{W})^2 = \sum_{i=1}^n (W_i - \bar{W})^2 = \sum_{i=1}^n (X_i - \bar{X})^2.$$

Now by Eq. (14) of [Chapter 1](#), the joint pdf of  $(Y_1, \dots, Y_n)$  is

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}) &= f_{\mathbf{W}}(\mathbf{A}^{-1}\mathbf{y}) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{A}^{-1}\mathbf{y})^\top(\mathbf{A}^{-1}\mathbf{y})\right] \\ &= \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left[-\frac{1}{2\sigma^2}\mathbf{y}^\top \mathbf{y}\right], \end{aligned}$$

that is,  $Y_1, \dots, Y_n$  are iid  $N(0, \sigma^2)$ . Thus

- (i)  $\sqrt{n}(\bar{X} - \mu) = \mathbf{a}_1^\top \mathbf{W} = Y_1$  is  $N(0, \sigma^2)$ , that is,  $\bar{X}$  is  $N(\mu, \sigma^2/n)$ ,
- (ii)  $\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (W_i - \bar{W})^2 = \sum_{i=2}^n Y_i^2$  is  $\sigma^2 \chi_{n-1}^2$ , that is,  $(n-1)s^2/\sigma^2$  is  $\chi_{n-1}^2$ ,
- (iii)  $\bar{X} = \mu + n^{-1/2}Y_1$  and  $s^2 = (n-1) \sum_{i=2}^n Y_i^2$  are independent.

□

## 2.2.8 $t$ and $F$ Distributions

**Definition 2.2.7.** A rv  $t_k = Z/\sqrt{W/k}$ , where  $Z \sim N(0, 1)$  and  $W \sim \chi_k^2$  are independent, is called a  $t$  rv with  $k$  df. It follows from [Proposition 2.2.5](#) that

$$\frac{\sqrt{n}(\bar{X} - \mu)}{s} = \frac{\sqrt{n}(\bar{X} - \mu)/\sigma}{\sqrt{\{(n-1)s^2/\sigma^2\}/(n-1)}}$$

is a  $t$  rv with  $(n-1)$  df.

This rv is commonly referred to as Student's  $t$  ("Student," [5]).

**Proposition 2.2.6.** If  $Z \sim N(0, 1)$  and  $W \sim \chi^2_k$  are independent, then  $T \sim t_k = Z/\sqrt{W/k}$  has pdf

$$f_T(t) = \frac{\Gamma((k+1)/2)}{\Gamma(1/2)\Gamma(k/2)} \frac{1}{\sqrt{k}} (1+t^2/k)^{-(k+1)/2}, \quad -\infty < t < \infty.$$

A distribution with this pdf is called the *t* distribution with  $k$  df.

*Proof.* The cdf of  $T \sim t_k$  is

$$\begin{aligned} F_T(t) &= P[Z/\sqrt{W/k} \leq t] = \int_{w=0}^{\infty} \int_{z=-\infty}^{t\sqrt{w/k}} f_Z(z)f_W(w) dz dw, \text{ and} \\ f_T(t) &= \int_0^{\infty} f_W(w) \sqrt{w/k} f_Z(t\sqrt{w/k}) dw \text{ (by Eq. (13) in Chapter 1)} \\ &= \int_0^{\infty} \left( \frac{1}{\Gamma(k/2)} 2^{-k/2} e^{-w/2} w^{k/2-1} \right) \sqrt{\frac{w}{k}} \left( \frac{1}{\sqrt{2\pi}} e^{-(1/2)t^2 w/k} \right) dw \\ &= \frac{(1/2)^{(k+1)/2}}{\Gamma(1/2)\Gamma(k/2)} \frac{1}{\sqrt{k}} \int_0^{\infty} e^{-(1/2)(1+t^2/k)w} w^{(k+1)/2-1} dw \\ &= \frac{\Gamma((k+1)/2)}{\Gamma(1/2)\Gamma(k/2)} \frac{1}{\sqrt{k}} (1+t^2/k)^{-(k+1)/2}, \quad -\infty < t < \infty, \end{aligned}$$

recognizing the last integral as  $\Gamma((k+1)/2)$  by appropriate normalization.  $\square$

Using Stirling's approximation,  $\Gamma(p) \approx \sqrt{2\pi}e^{-(p-1)}(p-1)^{p-1/2}$  for large  $p$ , we see that the pdf  $f_T(t)$  converges to the pdf of  $N(0, 1)$  at each  $t$ , as  $k \rightarrow \infty$ .

Another important rv is

$$F_{k_1, k_2} = \frac{\chi_{k_1}^2 / k_1}{\chi_{k_2}^2 / k_2},$$

where  $\chi_{k_1}^2$  and  $\chi_{k_2}^2$  are independent  $\chi^2$  rv's with  $k_1$  df and  $k_2$  df, respectively.

To find the pdf of  $F_{k_1, k_2}$ , we first use Eq. (13) and then Eq. (10) in Chapter 1 for the pdf of the ratio of two independent rv's and then for the pdf of the multiple of an rv. In this way, we arrive at

$$f_{F_{k_1, k_2}}(w) = \frac{\Gamma((k_1+k_2)/2)}{\Gamma(k_1/2)\Gamma(k_2/2)} \left( \frac{k_1}{k_2} \right)^{k_1/2} \frac{w^{k_1/2-1}}{[1+(k_1/k_2)w]^{(k_1+k_2)/2}}, \quad w > 0,$$

the verification of which is left as an exercise.

## 2.2.9 Noncentral $\chi^2$ and $F$ Distributions

In Section 2.2.5 we introduced the noncentral  $\chi^2$  rv with 1 df and noncentrality parameter  $\mu^2/2$  as the rv  $X^2$  where  $X \sim N(\mu, 1)$ . The pdf of this rv was expressed as  $E[f_{\chi_{1+2N}^2}(w)]$  where  $N$  is a  $Poi(\mu^2/2)$  rv. We now extend this by considering the distribution of  $\sum_{i=1}^k X_i^2$  where  $X_1, \dots, X_k$  are independent  $N(\mu_i, 1)$ .

As in the proof of the [Proposition 2.2.5](#), transform  $\mathbf{Y} = \mathbf{AX}$  where  $\mathbf{A} = ((a_{ij}))$  with  $\mathbf{a}_i^\top = (a_{i1}, \dots, a_{ik})$  satisfying

$$\mathbf{a}_1^\top = \frac{1}{\|\boldsymbol{\mu}\|}(\mu_1, \dots, \mu_k) \text{ and } \mathbf{a}_i^\top \mathbf{a}_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases},$$

where  $\|\boldsymbol{\mu}\|^2 = \sum_{i=1}^k \mu_i^2$ . Then  $\mathbf{A}$  is orthonormal and  $\sum_{i=1}^k X_i^2 = \sum_{i=1}^k Y_i^2$ . Moreover, since  $Y_1 = \mathbf{a}_1^\top \mathbf{X} = \sum_{i=1}^k \mu_i X_i / \|\boldsymbol{\mu}\|$ , we have

$$\begin{aligned} (\mathbf{A}^{-1} \mathbf{Y} - \boldsymbol{\mu})^\top (\mathbf{A}^{-1} \mathbf{Y} - \boldsymbol{\mu}) &= (\mathbf{X} - \boldsymbol{\mu})^\top (\mathbf{X} - \boldsymbol{\mu}) = \sum_{i=1}^k (X_i - \mu_i)^2 \\ &= \sum_{i=1}^k X_i^2 - 2 \sum_{i=1}^k \mu_i X_i + \sum_{i=1}^k \mu_i^2 = \sum_{i=1}^k Y_i^2 - 2\|\boldsymbol{\mu}\| Y_1 + \|\boldsymbol{\mu}\|^2 \\ &= (Y_1 - \|\boldsymbol{\mu}\|)^2 + \sum_{i=2}^k Y_i^2, \end{aligned}$$

and the joint pdf of  $(Y_1, \dots, Y_k)$  is

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(\mathbf{A}^{-1} \mathbf{y}) &= \frac{1}{(\sqrt{2\pi})^n} \exp\left[-\frac{1}{2}(\mathbf{A}^{-1} \mathbf{Y} - \boldsymbol{\mu})^\top (\mathbf{A}^{-1} \mathbf{Y} - \boldsymbol{\mu})\right] \\ &= \frac{1}{(\sqrt{2\pi})^n} \exp\left[-\frac{1}{2}(y_1 - \|\boldsymbol{\mu}\|)^2 - \frac{1}{2} \sum_{i=2}^k y_i^2\right]. \end{aligned}$$

It follows that  $Y_1 \sim N(\|\boldsymbol{\mu}\|, 1)$  and  $\sum_{i=2}^k Y_i^2 \sim \chi_{k-1}^2$ , which are independent. Now

$$\sum_{i=1}^k X_i^2 = Y_1^2 + \sum_{i=2}^k Y_i^2 \sim \chi_{1+2N}^2 + \chi_{k-1}^2,$$

where the two  $\chi^2$  rv's are independent. Using additive property of independent  $\chi^2$  rv's, we see that  $W = \sum_{i=1}^k X_i^2 = \chi_{k+2N}^2$  has the pdf

$$f_W(w) = \sum_{r=0}^{\infty} f_N(r, \|\boldsymbol{\mu}\|^2/2) f_{\chi_{k+2r}^2}(w) = E\left[f_{\chi_{k+2N}^2}(w)\right],$$

$N \sim Poi(\|\boldsymbol{\mu}\|^2/2)$ . We denote this by writing  $W \sim \chi_k^2(\|\boldsymbol{\mu}\|^2/2)$ , which is a  $\chi^2$  rv with  $k$  df and noncentrality parameter  $\|\boldsymbol{\mu}\|^2/2$ .

A rv of the form  $V = \frac{W_1/m}{W_2/n}$ , where  $W_1 \sim \chi_m^2(\delta^2)$  and  $W_2 \sim \chi_n^2$  are independent, is said to be distributed as  $F$  with numerator df  $m$ , denominator df  $n$  and noncentrality parameter  $\delta^2$ , denoted by  $F_{m,n}(\delta^2)$ .

Since  $W_1 \sim \chi^2_{m+2N}$  with  $N \sim Poi(\delta^2/2)$ , we can write  $V = F_{m+2N, n}$ , where  $F_{m+2r, n}$  is  $F$  with numerator df  $m+2r$  and denominator df  $n$ . In particular, if  $X_1, \dots, X_n$  are independent  $N(\mu, \sigma^2)$ , then

$$T = \frac{\sqrt{n}\bar{X}}{s} = \frac{\sqrt{n}\bar{X}/\sigma}{\sqrt{\{(n-1)s^2/\sigma^2\}/(n-1)}}$$

is distributed as

$$\frac{N(\sqrt{n}\mu/\sigma, 1)}{\sqrt{\chi^2_{n-1}/(n-1)}},$$

the numerator and denominator being independent (by [Proposition 2.2.5](#)), so that

$$T^2 \sim \frac{\chi^2_1(n\mu^2/\sigma^2)}{\chi^2_{n-1}/(n-1)} = F_{1, n-1}\left(\frac{n\mu^2}{\sigma^2}\right).$$

The noncentral  $\chi^2$  and  $F$  distributions are used in evaluating “power properties” of tests of significance of various hypotheses about means and variances of normal populations and in tests of hypotheses in Linear Models, as will be seen later. The noncentrality parameter is a measure of departure from “null hypothesis,” and for fixed  $c, k, m, n$ ,  $P[\chi^2_k(\delta^2) \geq c]$  and  $P[F_{m, n}(\delta^2) \geq c]$  are increasing functions of  $\delta^2$ . Proofs of these monotonicity properties are left as exercises.

### 2.2.10 Cauchy Distribution

**Definition 2.2.8.** A rv  $X$  is distributed as *Cauchy*( $\theta$ ) if  $X$  has pdf

$$f_X(x) = \frac{1}{\pi} \frac{1}{1 + (x - \theta)^2}, \quad -\infty < x < \infty.$$

We first verify that  $f_X$  is indeed a pdf. Since  $f_X(x) > 0$  for all  $x$ , we only need to check

$$\begin{aligned} \int_{-\infty}^{\infty} f_X(x) dx &= \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{1}{1 + (x - \theta)^2} dx = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{1}{1 + y^2} dy \\ &= \frac{2}{\pi} \int_0^{\infty} \frac{1}{1 + y^2} dy = \frac{2}{\pi} \int_0^{\pi/2} \frac{1 + \tan^2 \varphi}{1 + \tan^2 \varphi} d\varphi \\ &= \frac{2}{\pi} \int_0^{\pi/2} d\varphi = 1, \end{aligned}$$

using the transformation  $\varphi = \tan^{-1} y$ , that is,  $y = \tan \varphi$ . To find the mean of  $X$ , note that  $X = Y + \theta$  where  $Y = X - \theta$  has pdf

$$f_Y(y) = \frac{1}{\pi} \frac{1}{1 + y^2}, \quad -\infty < y < \infty.$$

Now

$$\begin{aligned} \mathbb{E}[|Y|] &= \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{|y|}{1+y^2} dy = \frac{2}{\pi} \int_0^{\infty} \frac{y}{1+y^2} dy \\ &= \frac{1}{\pi} \int_1^{\infty} \frac{1}{t} dt \quad (\text{with } t = 1+y^2), \end{aligned}$$

which diverges. Hence  $\mathbb{E}[Y]$  does not exist, and therefore,  $\mathbb{E}[X]$  does not exist.

### 2.2.11 Multivariate Normal Distribution

**Definition 2.2.9.** If  $X^T = (X_1, \dots, X_p)$  is a  $p$ -dim rv which can be written as  $X = \mu + BZ$ , where  $\mu$  is a  $p$ -dim vector in  $\mathbb{R}^p$ ,  $B$  is a  $p \times k$  matrix and  $Z$  is a  $k$ -dim vector of iid  $N(0, 1)$  rv's, then  $X$  is said to follow a  $p$ -variate normal distribution.

Since  $\mathbb{E}[X] = \mu$  and  $\text{Cov}[X, X] = \mathbb{E}[(X - \mu)(X - \mu)^T] = BB^T := \Sigma$ , we call  $X$  a  $p$ -variate normal rv with mean vector  $\mu$  and covariance matrix  $\Sigma$ , and write  $X \sim N_p(\mu, \Sigma)$ .

We now list the important properties of  $p$ -variate normal distribution in the following propositions.

**Proposition 2.2.7.** If  $X \sim N_p(\mu, \Sigma)$  and  $Y = c + AX$ , where  $c$  is in  $\mathbb{R}^r$  and  $A$  is a  $r \times p$  matrix, then  $Y \sim N_r(c + A\mu, A\Sigma A^T)$ .

It follows that the vector formed by any subset of  $r$  coordinates,  $1 \leq r \leq p-1$ , follows an  $r$ -variate normal distribution with appropriate mean vector and covariance matrix. In particular, each  $X_i \sim N(\mu_i, \sigma_{ii})$ .

*Proof.* Since  $X = \mu + BZ$ , the transform

$$Y = c + A(\mu + BZ) = (c + A\mu) + (AB)Z,$$

and  $(AB)(AB)^T = A(BB^T)A^T = A\Sigma A^T$ . Hence,  $Y \sim N_r(c + A\mu, A\Sigma A^T)$ .  $\square$

**Proposition 2.2.8.** If  $p \leq k$  and  $\text{rank}(B) = p$ , then  $\Sigma = BB^T$  is positive definite, and  $X = \mu + BZ$  has pdf

$$f_X(x; \mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right], \quad x \in \mathbb{R}^p.$$

Conversely, if  $X$  has pdf  $f_X(x; \mu, \Sigma)$  given above, for some  $\mu$  and positive definite  $\Sigma$ , then  $X$  can be expressed as  $X = \mu + BZ$ ,  $Z^T = (Z_1, \dots, Z_p)$  with  $Z_1, \dots, Z_p$  iid as  $N(0, 1)$ .

*Proof.* Since the rank of  $B$  is  $p \leq k$ ,  $B^T a \neq \mathbf{0}$  for any  $a \neq \mathbf{0}$  in  $\mathbb{R}^p$ , so that  $\Sigma = BB^T$  is positive definite. Next augment  $B$  and  $\mu$  with a  $(k-p) \times k$  matrix  $C$  and a  $k-p$  vector of zeros, respectively, so that

$$B^* = \begin{bmatrix} B \\ C \end{bmatrix} \text{ and } \mu^* = \begin{bmatrix} \mu \\ \mathbf{0} \end{bmatrix},$$

where the rows of  $C$  are of unit length, mutually orthogonal, and orthogonal to the rows of  $B$ . Note that  $B^*$  is a square matrix of order  $k$  and  $\mu^*$  is a  $k$ -dim vector. Then letting  $Y = CZ$ , we have

$$\mathbf{X}^* = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} + \mathbf{BZ} \\ \mathbf{CZ} \end{bmatrix}$$

as a one-to-one transformation from  $\mathbf{Z}$ . We now find the pdf of  $\mathbf{X}^*$  from the pdf of  $\mathbf{Z}$  and integrate out  $\mathbf{Y}$ . The details are left as an exercise.

To prove the converse, find  $\mathbf{B}$  so that  $\boldsymbol{\Sigma} = \mathbf{BB}^\top$  ( $\mathbf{B}$  is not unique) and transform  $\mathbf{Z} = \mathbf{B}^{-1}(\mathbf{X} - \boldsymbol{\mu})$ . Then  $\mathbf{Z}$  has the desired property.  $\square$

*Remark 2.2.1.* If  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are  $p$ - and  $q$ -dimensional rv's, then independence of  $\mathbf{X}_1$  and  $\mathbf{X}_2$  implies  $\boldsymbol{\Sigma}_{12} = \text{Cov}[\mathbf{X}_1, \mathbf{X}_2] = 0$ . However, the converse is not true in general. In the case of normal distributions the converse is also true. More precisely, if  $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \sim N_{p+q}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are independent iff  $\boldsymbol{\Sigma}_{12} = \text{Cov}[\mathbf{X}_1, \mathbf{X}_2] = 0$ . The proof can be easily seen when  $\boldsymbol{\Sigma}$  is positive definite. If  $\boldsymbol{\Sigma}_{12} = 0$ , then  $\boldsymbol{\Sigma}$  is a block diagonal matrix of the form  $\begin{bmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0}^\top & \boldsymbol{\Sigma}_{22} \end{bmatrix}$  and its inverse is also block diagonal. This results in the joint pdf of  $\mathbf{X}_1$  and  $\mathbf{X}_2$  being the product of the pdf's of  $\mathbf{X}_1$  and  $\mathbf{X}_2$ .

**Proposition 2.2.9.** Suppose  $\mathbf{X} \sim N_{p+q}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\Sigma}$  is positive definite. Let

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix},$$

where  $\mathbf{X}_1$ ,  $\boldsymbol{\mu}_1$  and  $\mathbf{X}_2$ ,  $\boldsymbol{\mu}_2$  denote the first  $p$  and last  $q$  coordinates of  $\mathbf{X}$  and  $\boldsymbol{\mu}$ , respectively, and  $\boldsymbol{\Sigma}_{11} = \text{Cov}[\mathbf{X}_1, \mathbf{X}_1]$ ,  $\boldsymbol{\Sigma}_{22} = \text{Cov}[\mathbf{X}_2, \mathbf{X}_2]$ ,  $\boldsymbol{\Sigma}_{21}^\top = \boldsymbol{\Sigma}_{12} = \text{Cov}[\mathbf{X}_1, \mathbf{X}_2]$ , where  $\text{Cov}[\mathbf{X}_i, \mathbf{X}_j] = E[(\mathbf{X}_i - \boldsymbol{\mu}_i)(\mathbf{X}_j - \boldsymbol{\mu}_j)^\top]$ . Then the conditional distribution of  $\mathbf{X}_2$  given  $\mathbf{X}_1 = \mathbf{x}_1$  is  $N_q(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})$ .

*Proof.* Write  $f_{\mathbf{X}_2|\mathbf{X}_1}(\mathbf{x}_2|\mathbf{x}_1) = f_{\mathbf{X}_1, \mathbf{X}_2}(\mathbf{x}_1, \mathbf{x}_2)/f_{\mathbf{X}_1}(\mathbf{x}_1)$ , where

$$\begin{aligned} f_{\mathbf{X}_1, \mathbf{X}_2}(\mathbf{x}_1, \mathbf{x}_2) &= \frac{1}{(2\pi)^{(p+q)/2} \left| \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right|^{1/2}} \\ &\times \exp \left\{ -\frac{1}{2} [(\mathbf{x}_1 - \boldsymbol{\mu}_1)^\top, (\mathbf{x}_2 - \boldsymbol{\mu}_2)^\top] \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x}_1 - \boldsymbol{\mu}_1 \\ \mathbf{x}_2 - \boldsymbol{\mu}_2 \end{bmatrix} \right\} \end{aligned}$$

and

$$f_{\mathbf{X}_1}(\mathbf{x}_1) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_{11}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_1 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \right].$$

To calculate the ratio of the last two terms, the main task is to find the inverse and the determinant of the partitioned matrix involved in  $f_{\mathbf{X}_1, \mathbf{X}_2}(\mathbf{x}_1, \mathbf{x}_2)$ . For this we need the following argument. Let

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \text{ with } \mathbf{A}_{21} = \mathbf{A}_{12}^\top \text{ and } \mathbf{B} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{A}_{21}^\top \mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix}.$$

Then

$$\mathbf{D} = \mathbf{BAB}^\top = \begin{bmatrix} A_{11} & \mathbf{0} \\ \mathbf{0} & A_{22} - A_{12}^\top A_{11}^{-1} A_{12} \end{bmatrix},$$

so  $\mathbf{D}^{-1} = \begin{bmatrix} A_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & (A_{22} - A_{12}^\top A_{11}^{-1} A_{12})^{-1} \end{bmatrix}$ . This implies  $\mathbf{A} = \mathbf{B}^{-1} \mathbf{D} (\mathbf{B}^\top)^{-1}$  and therefore, denoting  $(A_{22} - A_{12}^\top A_{11}^{-1} A_{12})^{-1}$  by  $\mathbf{H}$ ,

$$\mathbf{A}^{-1} = \mathbf{B}^\top \mathbf{D}^{-1} \mathbf{B} = \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1} A_{12} \mathbf{H} A_{12}^\top A_{11}^{-1} & -A_{11}^{-1} A_{12} \mathbf{H} \\ -\mathbf{H} A_{12}^\top A_{11}^{-1} & \mathbf{H} \end{bmatrix}.$$

Finally,

$$|\mathbf{A}| = |\mathbf{B}|^{-1} |\mathbf{D}| |\mathbf{B}^\top|^{-1} = |\mathbf{D}| = |A_{11}| |A_{22} - A_{12}^\top A_{11}^{-1} A_{12}|,$$

because  $|\mathbf{B}| = |\mathbf{B}^\top| = 1$ . Using these results on  $\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}^{-1}$  and  $\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ , we obtain after some algebraic simplification,

$$\begin{aligned} f_{\mathbf{X}_2|\mathbf{X}_1}(\mathbf{x}_2|\mathbf{x}_1) &= f_{\mathbf{X}_1, \mathbf{X}_2}(\mathbf{x}_1, \mathbf{x}_2)/f_{\mathbf{X}_1}(\mathbf{x}_1) \\ &= \frac{1}{(2\pi)^{q/2} |\Sigma_{22,1}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x}_2 - \mu_{2,1})^\top \Sigma_{22,1}^{-1} (\mathbf{x}_2 - \mu_{2,1})\right], \end{aligned}$$

where  $\mu_{2,1} = \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (\mathbf{x}_1 - \mu_1)$  and  $\Sigma_{22,1} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ . The expression displayed above is the pdf of  $N_q(\mu_{2,1}, \Sigma_{22,1})$ .  $\square$

*Remark 2.2.2. Special Cases.*

(i) For  $q = 1$ , write  $\sigma_{2(1)}^\top = (\sigma_{1,p+1}, \dots, \sigma_{p,p+1})$  where  $\sigma_{ij} = \text{Cov}[X_i, X_j]$ . Then

$$\begin{aligned} \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} &= \sigma_{p+1}^2 - \sigma_{2(1)}^\top \Sigma_{11}^{-1} \sigma_{2(1)} := \sigma_{(p+1)\cdot(12\dots p)}^2, \text{ and} \\ \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (\mathbf{x}_1 - \mu_1) &= (\mu_{p+1} - \sigma_{2(1)}^\top \Sigma_{11}^{-1} \mu_1) + \sigma_{2(1)}^\top \Sigma_{11}^{-1} \mathbf{x}_1 \\ &:= \alpha + \beta_1 x_1 + \dots + \beta_p x_p, \end{aligned}$$

where  $\beta^\top = (\beta_1, \dots, \beta_p) = \sigma_{2(1)}^\top \Sigma_{11}^{-1}$  and  $\alpha = \mu_{p+1} - (\beta_1 \mu_1 + \dots + \beta_p \mu_p)$ . (Remember: In the above,  $\mathbf{x}_1^\top = (x_1, \dots, x_p)$  and  $\mu_1^\top = (\mu_1, \dots, \mu_p)$ .) Thus the conditional distribution of  $X_{p+1}$  given  $(X_1, \dots, X_p) = (\mathbf{x}_1, \dots, \mathbf{x}_p)$  is

$N(\alpha + \beta_1 x_1 + \dots + \beta_p x_p, \sigma_{(p+1)\cdot(12\dots p)}^2)$  where the conditional mean  $\alpha + \beta_1 x_1 + \dots + \beta_p x_p$  and the conditional variance  $\sigma_{(p+1)\cdot(12\dots p)}^2$  are given by the above formulas.

(ii) For  $p = q = 1$ ,  $\Sigma = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$  and therefore,

$$|\Sigma|^{1/2} = \sigma_1 \sigma_2 \sqrt{1 - \rho^2} \text{ and } \Sigma^{-1} = (1 - \rho^2)^{-1} \begin{bmatrix} 1/\sigma_1^2 & -\rho/(\sigma_1 \sigma_2) \\ -\rho/(\sigma_1 \sigma_2) & 1/\sigma_2^2 \end{bmatrix}.$$

Thus for  $-1 < \rho < 1$ ,

$$f_{X_1, X_2}(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x_1 - \mu_1}{\sigma_1}\right)\left(\frac{x_2 - \mu_2}{\sigma_2}\right) + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2 \right\} \right],$$

which is the pdf of bivariate normal distribution  $N_2(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$  with mean vector  $(\mu_1, \mu_2)^\top$  and covariance matrix  $\begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$ . The marginal distribution of  $X_i$  is  $N(\mu_i, \sigma_i^2)$  and the conditional distribution of  $X_2$  given  $X_1 = x_1$  is  $N(\alpha + \beta x_1, \sigma_{2,1}^2)$ , where  $\beta = \rho\sigma_2/\sigma_1$ ,  $\alpha = \mu_2 - \beta\mu_1$ , and  $\sigma_{2,1}^2 = \sigma_2^2(1 - \rho^2)$ . Conversely, if  $X_1 \sim N(\mu_1, \sigma_1^2)$  and the conditional distribution of  $X_2$  given  $X_1$  is  $N(\alpha + \beta X_1, \tau^2)$ , then  $(X_1, X_2) \sim N_2(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$ , where  $\mu_2 = \alpha + \beta\mu_1$ ,  $\sigma_2^2 = \tau^2 + \beta^2\sigma_1^2$ , and  $\rho = \beta\sigma_1/\sqrt{\tau^2 + \beta^2\sigma_1^2}$ . Moreover, if  $(X_1, X_2) \sim N_2(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$  and  $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$  is nonsingular, then  $Y_1 = a_{11}X_1 + a_{12}X_2$  and  $Y_2 = a_{21}X_1 + a_{22}X_2$  are jointly bivariate normal with appropriate parameters. In particular  $X_1 + X_2$  and  $X_1 - X_2$  are independent iff  $\sigma_1 = \sigma_2$ .

**Proposition 2.2.10.** Suppose  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with positive definite  $\boldsymbol{\Sigma}$ . Then  $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \sim \chi_p^2$ .

*Proof.* Find  $\mathbf{B}$  such that  $\mathbf{B}\mathbf{B}^\top = \boldsymbol{\Sigma}$  and write  $\mathbf{X} = \boldsymbol{\mu} + \mathbf{BZ}$  with  $\mathbf{Z} \sim N_p(\mathbf{0}, \mathbf{I})$ , as in the proof of the converse part of [Proposition 2.2.8](#). Then

$$(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) = \mathbf{Z}^\top \mathbf{B}^\top (\mathbf{B}\mathbf{B}^\top)^{-1} \mathbf{BZ} = \mathbf{Z}^\top \mathbf{Z} = \sum_{i=1}^p Z_i^2 \sim \chi_p^2.$$

□

## 2.2.12 Exponential Family of Distributions

A family of pdf's or pmf's is said to be an exponential family if it is of the form

$$f(x, \boldsymbol{\theta}) = c(\boldsymbol{\theta}) \exp\left[ \sum_{j=1}^k Q_j(\boldsymbol{\theta}) T_j(x) \right] r(x), \quad \boldsymbol{\theta} \in \Theta,$$

where  $r(x) > 0$ ,  $c(\boldsymbol{\theta}) > 0$ , and each  $T_j(x)$  is a real-valued function of  $x$ , while each  $Q_j(\boldsymbol{\theta})$  is a real-valued function of  $\boldsymbol{\theta}$ .

**Note.** The set  $\{x: f(x, \boldsymbol{\theta}) > 0\} = \{x: r(x) > 0\}$  does not depend on  $\boldsymbol{\theta}$ .

Many distributions including most of those discussed above belong to this family.

**Example 2.2.2.** The pmf of the Binomial distribution  $\text{Bin}(n, p)$  is

$$\begin{aligned} f(x, p) &= \binom{n}{x} \exp[x \log(p) + (n-x) \log(1-p)] \\ &= \binom{n}{x} (1-p)^n \exp\left[x \log\left(\frac{p}{1-p}\right)\right], \end{aligned}$$

$x = 0, 1, \dots, n$ . Thus  $f(x, p)$  is of the desired form with  $\theta = p$ ,  $c(\theta) = (1-\theta)^n$ ,  $r(x) = \binom{n}{x}$ ,  $Q(\theta) = \log(\theta/(1-\theta))$ , and  $T(x) = x$ .

**Example 2.2.3.** The pmf of  $\text{Poi}(\lambda)$  is

$$f(x, \lambda) = e^{-\lambda} \frac{\lambda^x}{x!} = e^{-\lambda} \frac{1}{x!} \exp(x \log \lambda), \quad x = 0, 1, \dots$$

This pmf is of the desired form with  $\theta = \lambda$ ,  $r(x) = 1/x!$ ,  $c(\theta) = e^{-\theta}$ ,  $Q(\theta) = \log \theta$ , and  $T(x) = x$ .

**Example 2.2.4.** The pdf of  $N(\mu, \sigma^2)$  is

$$\begin{aligned} f(x; \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right] \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\mu^2/(2\sigma^2)}\right) \exp\left[\frac{\mu}{\sigma^2}x + \left(-\frac{x^2}{2\sigma^2}\right)\right]. \end{aligned}$$

This pdf is of the desired form with  $\theta = (\mu, \sigma^2)$ ,  $r(x) = 1$ ,  $c(\theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\mu^2/(2\sigma^2)}$ ,  $Q_1(\theta) = \mu/\sigma^2$ ,  $T_1(x) = x$ ,  $Q_2(\theta) = -1/(2\sigma^2)$ , and  $T_2(x) = x^2$ .

The pdf of  $\text{Unif}(0, \theta)$  is  $f(x, \theta) = I_{(0,\theta)}(x)$ . Here, the set  $\{x: f(x, \theta) > 0\} = (0, \theta)$  depends on  $\theta$ . Therefore, this family of pdf's  $\{f(x, \theta), \theta > 0\}$  is not an exponential family.

*Natural Parameters.* In the form of  $f(x; \theta)$  given in the definition of an exponential family, if we reparameterize by letting

$$\theta_j = Q_j(\theta), \quad j = 1, \dots, k, \text{ and let } c(\theta) = c(\theta_1, \dots, \theta_k),$$

then the  $k$ -parameter exponential family of pdf's or pmf's is represented as

$$f(x; \theta) = c(\theta) \exp\left[\sum_{j=1}^k \theta_j T_j(x)\right] r(x) = c(\theta) \exp[\langle \theta, T(x) \rangle] r(x),$$

where the parameters  $\theta_1, \dots, \theta_k$  are called *natural parameters* and  $\Theta^*$  is called the natural parameter space defined as

$$\Theta^* = \left\{ \theta = (\theta_1, \dots, \theta_k): \int_{-\infty}^{\infty} \exp\left[\sum_{j=1}^k \theta_j T_j(x)\right] r(x) dx < \infty \right\},$$

where  $\mathcal{X} = \mathbb{R}$  (and analogously for more general  $\mathcal{X}$ ). Also, the constant  $c(\theta)$  is given by

$$c(\theta) = \left\{ \int_{-\infty}^{\infty} \exp\left[\sum_{j=1}^k \theta_j T_j(x)\right] r(x) dx \right\}^{-1}.$$

With the reparameterization  $\theta_1 = \mu/\sigma^2$ ,  $\theta_2 = -1/(2\sigma^2)$ , the pdf's of the family of normal distributions can be written as

$$f(x; \theta_1, \theta_2) = c(\theta) \exp[\theta_1 x + \theta_2 x^2].$$

**Proposition 2.2.11.** *The natural parameter space  $\Theta^*$  of a  $k$ -parameter exponential family of distributions is a convex set in  $\mathbb{R}^k$ ; that is, if  $\theta$  and  $\theta'$  are in  $\Theta^*$  then  $\alpha\theta + (1 - \alpha)\theta'^*$  for all  $0 \leq \alpha \leq 1$ .*

*Proof.* We shall use Hölder's Inequality: if  $f$  and  $g$  are real-valued functions on  $\mathbb{R}^n$  for which  $\int |f|^p$  and  $\int |g|^q$  are finite for  $p > 1$ ,  $q > 1$  with  $1/p + 1/q = 1$ , then  $\int |fg| \leq (\int |f|^p)^{1/p} (\int |g|^q)^{1/q}$  (see Section A.2).

To prove the proposition, let  $p = 1/\alpha$ ,  $q = 1/(1 - \alpha)$  for  $0 < \alpha < 1$ , so that  $p$  and  $q$  are larger than 1 with  $1/p + 1/q = 1$ . Then by Hölder's inequality,

$$\int \exp[\langle \theta, T(x) \rangle] r(x) dx < \infty \text{ and } \int \exp[\langle \theta', T(x) \rangle] r(x) dx < \infty$$

and hence

$$\begin{aligned} & \int \exp[\langle \alpha\theta + (1 - \alpha)\theta', T(x) \rangle] r(x) dx \\ &= \int \{\exp[\alpha\langle \theta, T(x) \rangle] r^\alpha(x)\} \{\exp[(1 - \alpha)\langle \theta', T(x) \rangle] r^{1-\alpha}(x)\} dx \\ &\leq \left[ \int \{\exp[\alpha\langle \theta, T(x) \rangle] r^\alpha(x)\}^{1/\alpha} dx \right]^\alpha \\ &\quad \times \left[ \int \{\exp[(1 - \alpha)\langle \theta', T(x) \rangle] r^{1-\alpha}(x)\}^{1/(1-\alpha)} dx \right]^{1-\alpha} \\ &= \left[ \int \exp[\langle \theta, T(x) \rangle] r(x) dx \right]^\alpha \left[ \int \exp[\langle \theta', T(x) \rangle] r(x) dx \right]^{1-\alpha} \\ &< \infty. \end{aligned}$$

□

## Exercises

- 2.1.** A multiple choice exam consists of 20 questions, each with four possible answers and carries four points. If a question is attempted, then the score is 4 for a correct answer and 0 otherwise, while a score of 1 is awarded for a question which is not attempted. A student knows the answers to eight questions, but has no idea about six questions and is 50% sure about the answers of the other six questions. It takes 50 points to pass the exam. Find the probability of passing for each of the following strategies and choose the best one.

- (a) Answer the 8 sure ones and guess at random the answers of the other 12.
- (b) Answer the 8 sure ones, choose the answers of the 6 that you are 50% sure about and guess at random the other 6.
- (c) Answer the 8 sure ones and do not answer the other 12.

- 2.2.** Samples from a large lot of manufactured items are being inspected to determine whether the proportion of defectives  $p = 0.1$  or less as claimed by the supplier. The following sampling plan is under consideration:

Take a sample of 15 items and accept the lot if all are good, reject the lot if 2 or more are defectives and take another sample of 15 items if there is 1 defective in the sample. Accept the lot if all 15 in the second sample are good and reject the lot otherwise.

Find the probabilities that this sampling plan will result in

- (a) rejection of a lot with  $p = 0.1$ ,
- (b) acceptance of a lot with  $p = 0.15$ .

- 2.3.** An insurance company writes automobile insurance policies in an area where 1 out of 200 drivers causes accidents in a year according to past record. If the company writes 500 policies, find

- (a) the probability that there will be no more than two claims in a year,
- (b) the expected amount of claims to be settled if the average claim is \$1500.

- 2.4.** On a certain segment of a freeway, accidents happen during the rush hour at the rate of two per hour, following the Poisson distribution. Find the probabilities of

- (a) at least one accident in an hour,
- (b) if there was one accident in an hour, the accident occurred during the first half hour.

- 2.5.** A radioactive source emits particles at a rate of 0.4/s. Suppose that the number of particles emitted per second is a Poisson rv.

- (a) Find the probability that two or more particles will be emitted in 3 s.
- (b) A counter registers an emitted particle with probability 0.75. What is the probability that two or more particles will be registered on the counter in 3 s?

- 2.6.** Three players A, B, and C will play a table tennis match. In each game, A defeats B with probability 0.6, C defeats B with probability 0.6, and A and C are evenly matched. First, A plays against B and the one who wins three games, plays against C. Then the one who wins three games wins the match.

- (a) Find the probability that A wins the first round and the probability that the first round is settled in four games.
- (b) Find the probabilities of A, B, and C winning the match and the probability that the entire match is settled in nine games.
- (c) Find the expected number of games to settle the match.

- 2.7.** Find the moment generating function (mgf) of each of the following rv's:

- (a)  $X \sim Poi(\lambda)$ , (b)  $Y \sim Geom(p)$ , (c)  $T \sim Exp(\theta)$ , (d)  $W \sim Gamma(\alpha, \beta)$ .

Also find the means and variances of these rv's from their mgf's.

- 2.8.** Show that the pdf of  $F_{k_1, k_2} = (\chi^2_{k_1}/k_1)/(\chi^2_{k_2}/k_2)$  where  $\chi^2_{k_1}$  and  $\chi^2_{k_2}$  are independent chi-square rv's with df's  $k_1$  and  $k_2$ , respectively, is as given in the text.

- 2.9.** As in Section 2.2.8, let  $\chi_k^2(\delta^2)$  and  $F_{m,n}(\delta^2)$  denote the noncentral  $\chi_k^2$  and  $F_{m,n}$  rv's with noncentrality parameter  $\delta^2$ . Show that for fixed  $c, k, m, n$ ,  $P[\chi_k^2(\delta^2) \geq c]$  and  $P[F_{m,n}(\delta^2) \geq c]$  are increasing functions of  $\delta^2$ .
- 2.10.** Verify the details of the proof of Proposition 2.2.8.
- 2.11.** Let  $X_1, \dots, X_n$  be independent  $N(0, \sigma^2)$  rv's. Then  $Y_n = (X_1^2 + \dots + X_n^2)^{1/2}$  is distributed as  $\sigma\sqrt{\chi_n^2}$ . Find the pdf of  $Y_n$ . [This distribution is called Raleigh distribution for  $n = 2$  and Maxwell distribution for  $n = 3$ .]
- 2.12.** Let  $X_1, X_2$  be independent  $N(0, \sigma^2)$  rv's. Find the pdf's of  $Y_1 = X_1/X_2$ ,  $Y_2 = X_1/|X_2|$ ,  $Y_3 = |X_1|/|X_2|$ .
- 2.13.** Let  $V_1, \dots, V_{n+1}$  be independent  $\text{Exp}(1)$  rv's and let  $S_k = V_1 + \dots + V_k$ ,  $k = 1, \dots, n+1$ . Let  $U_{(1)} < \dots < U_{(n)}$  be the order statistics in a random sample  $(U_1, \dots, U_n)$  from  $\text{Unif}(0, 1)$ .
- (a) Show that the joint distribution of  $(U_{(1)}, \dots, U_{(n)})$  is the same as that of  $(S_1/S_{n+1}, \dots, S_n/S_{n+1})$ .
  - (b) Use this result and Proposition 2.2.2 to find the pdf's of (i)  $U_{(k)}$  and (ii)  $U_{(l)} - U_{(k)}$  for  $1 \leq k < l \leq n$ .
- 2.14.** Let  $X_{(1)} < \dots < X_{(n)}$  be the order statistics in a random sample  $(X_1, \dots, X_n)$  from  $\text{Exp}(1)$ . Show that  $nX_{(1)}, (n-1)(X_{(2)} - X_{(1)}), (n-2)(X_{(3)} - X_{(2)}), \dots, (X_{(n)} - X_{(n-1)})$  are iid  $\text{Exp}(1)$ .
- 2.15.**
  - (a) Show that if  $X$  has a continuous strictly increasing cdf  $F$ , then  $F(X)$  is  $\text{Unif}(0, 1)$ .
  - (b) If  $T_1, \dots, T_n$  are iid  $\text{Exp}(\text{mean } \theta)$ , then  $2(T_1 + \dots + T_n)/\theta$  is  $\chi_{2n}^2$ .
  - (c) Let  $(Z_1, Z_2) = (R \cos \theta, R \sin \theta)$  define a one-to-one map between  $(Z_1, Z_2)$  and  $(R, \theta)$ . What is the joint distribution of  $(R, \theta)$  when  $Z_1, Z_2$  are independent  $\text{N}(0, 1)$ ?
  - (d) Let  $Z_1, Z_2$  be independent  $N(0, 1)$  and let

$$X_1 = a_1 + b_{11}Z_1 + b_{12}Z_2, \quad X_2 = a_2 + b_{21}Z_1 + b_{22}Z_2$$

Find the constants  $a_1, a_2, b_{11}, b_{12}, b_{21}, b_{22}$  so that  $(X_1, X_2)$  follows a bivariate normal distribution with mean vector  $(\mu_1, \mu_2)$  and covariance matrix

$$\begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}.$$

- 2.16.** Suppose that you have a computer program to generate  $\text{Unif}(0, 1)$  rv's  $U_1, U_2, \dots$  How would you use these rv's to generate independent rv's distributed as
- (a)  $T \sim \text{Exp}(1)$ ,
  - (b)  $V \sim \chi_{10}^2$ ,
  - (c)  $W \sim \text{Cauchy}(0, 1)$ ,
  - (d)  $Z \sim N(0, 1)$ ,
  - (e)  $(X_1, X_2) \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\mu} = (\mu_1, \mu_2)$  and  $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$

Use the results in Exercise 2.15.

- 2.17.** Suppose that  $(X_1, X_2, X_3)$  follows a 3-dim normal distribution with the mean vector  $(\mu_1, \mu_2, \mu_3)$  and a positive definite covariance matrix  $\Sigma = ((\sigma_{ij}))$ . Let  $X_{2.1} = X_2 - E[X_2|X_1]$ ,  $X_{3.1} = X_3 - E[X_3|X_1]$ , and  $X_{3.12} = X_3 - E[X_3|X_1, X_2]$ .

- (a) Show that  $X_1, X_{2.1}$  and  $X_{3.12}$  are mutually independent  
 (b) The partial correlation coefficient between  $X_2, X_3$  given  $X_1$  is

$$\rho_{23.1} = \text{Corr}[X_{2.1}, X_{3.1}] = \frac{\text{Cov}[X_{2.1}, X_{3.1}]}{\sqrt{\text{Var}[X_{2.1}]\text{Var}[X_{3.1}]}}$$

Express  $\rho_{23.1}$  in terms of  $\rho_{12}, \rho_{13}$  and  $\rho_{23}$  where  $\rho_{ij} = \text{Corr}[X_i, X_j]$ .

- (c) The multiple correlation of  $X_3$  on  $(X_1, X_2)$  is  $\rho_{3.12} = \text{Corr}[X_3, E(X_3|X_1, X_2)]$ . Show that  $\rho_{3.12} \geq 0$  and  $\text{Var}[X_{3.12}] = \sigma_{33}(1 - \rho_{3.12}^2)$ .  
 (d) Show that  $1 - \rho_{3.12}^2 = (1 - \rho_{13}^2)(1 - \rho_{23.1}^2)$ .

# Infinite Sequences of Random Variables and Their Convergence Properties

## 3.1 Introduction

Let  $T$  be a statistic based on the data consisting of rv's  $X_1, \dots, X_n$ . Most statistical methods (such as estimation and hypothesis testing discussed in subsequent chapters) use a statistic appropriate for the problem at hand. It is of interest to know how  $T$  behaves as we have more and more data in order to understand the behavior of the procedures based on  $T$  with large data sets. For a rigorous examination of this question, we need proper notation for  $T$  based on  $(X_1, \dots, X_n)$  as  $n \rightarrow \infty$  and the concept of how  $T$  behaves as  $n \rightarrow \infty$  should be made more precise.

## 3.2 Modes of Convergence

Let  $T$  be a statistic based on the data consisting of rv's  $X_1, \dots, X_n$ . Most statistical methods (such as estimation and hypothesis testing discussed in subsequent chapters) use a statistic appropriate for the problem at hand. It is of interest to know how  $T$  behaves as we have more and more data in order to understand the behavior of the procedures based on  $T$  with large data sets. For a rigorous examination of this question, we need proper notation for  $T$  based on  $(X_1, \dots, X_n)$  as  $n \rightarrow \infty$  and the concept of how  $T$  behaves as  $n \rightarrow \infty$  should be made more precise.

**Definition 3.2.1.** Let  $X_1, \dots, X_n$  be iid rv's with common cdf  $F \in \mathcal{F}$ , where  $\mathcal{F}$  is the class of all cdf's on  $\mathbb{R}$ . Then the function

$$F_n(x) = n^{-1} \sum_{i=1}^n I_{(-\infty, x]}(X_i)$$

is called the empirical distribution function (edf).

**Definition 3.2.2.** Let  $\mathbf{T}: \mathcal{F} \rightarrow \mathbb{R}^k$ . Then  $\mathbf{T}_n = (T_{n1}, \dots, T_{nk}) = \mathbf{T}(F_n)$  is the  $k$ -dim statistic  $\mathbf{T}$  based on  $(X_1, \dots, X_n)$ , that is,  $\mathbf{T}_n = \mathbf{T}(X_1, \dots, X_n)$ .

We now define various modes of convergence.

**Definition 3.2.3.** The various modes of convergence are defined as:

1. Convergence in Probability:  $X_n \xrightarrow{P} a$  if  $\lim_{n \rightarrow \infty} P[|X_n - a| > \varepsilon] = 0$  for all  $\varepsilon > 0$  and  $X_n \xrightarrow{P} X$  if  $|X_n - X| \xrightarrow{P} 0$ , denoted by  $X_n = a + o_P(1)$  and  $X_n = X + o_P(1)$ , respectively.
2. Convergence in Law:  $X_n \xrightarrow{\mathcal{L}} X$  if  $P[X_n \in A] \rightarrow P[X \in A]$  for all  $A$  for which  $P[X \in \partial A] = 0$  where  $\partial A$  is the boundary of  $A$ . For 1-dim rv's,  $F_{X_n}(x) \rightarrow F_X(x)$  at all continuity points of  $F_X$  is enough.
3. Convergence in Quadratic Mean:  $X_n \xrightarrow{q.m.} a$  if  $E[(X_n - a)^2] \rightarrow 0$  and  $X_n \xrightarrow{q.m.} X$  if  $E[(X_n - X)^2] \rightarrow 0$ . More generally,  $X_n \xrightarrow{r} X$  if  $E[|X_n - X|^r] \rightarrow 0$  for  $X_n, X$  such that  $E[|X_n|^r] < \infty$  and  $E[|X|^r] < \infty$ .
4. Bounded in Probability:  $\{X_n\}$  is bounded in probability, denoted by  $X_n = O_P(1)$ , if for any  $\varepsilon > 0$  there exists a positive constant  $M_\varepsilon$  such that  $P[|X_n| \leq M_\varepsilon] > 1 - \varepsilon$  for all  $n$ .
5. Almost Sure Convergence:  $X_n \xrightarrow{a.s.} a$  if  $P[\lim_{n \rightarrow \infty} X_n = a] = 1$ .
6. Uniform Integrability:  $\{X_n\}$  is uniformly integrable if

$$\lim_{k \rightarrow \infty} \overline{\lim}_{n \rightarrow \infty} E[|X_n| I_{[k, \infty)}(|X_n|)] = 0.$$

Compare the stochastic order relations  $o_P$  and  $O_P$  with the usual order relations  $o$  and  $O$  for sequence of real numbers:

1. If  $\lim_{n \rightarrow \infty} x_n = 0$ , then  $x_n = o(1)$ .
2. If  $\{x_n\}$  is bounded, then  $x_n = O(1)$ .

*Extensions of the above definitions.* Let  $\{r_n\}$  be a sequence of real numbers. Then  $X_n = o_P(r_n)$  iff  $X_n/r_n = o_P(1)$ ,  $X_n = O_P(r_n)$  iff  $X_n/r_n = O_P(1)$ ,  $x_n = o(r_n)$  iff  $x_n/r_n = o(1)$ ,  $x_n = O(r_n)$  iff  $X_n/r_n = O(1)$ . For example,  $X_n = \theta + O_P(1/\sqrt{n})$  iff  $\sqrt{n}(X_n - \theta) = O_P(1)$ .

*Convergence Properties of Sample Means of iid Random Variables and Random Vectors.* Suppose  $\{X_n\}$  is a sequence of iid rv's and  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ .

**Theorem 3.2.1** (Weak Law of Large Numbers (WLLN, Khinchine)). *If  $E[X_1] = \mu$  exists, then  $\bar{X}_n \xrightarrow{P} \mu$ .*

**Theorem 3.2.2** (Central Limit Theorem (CLT, Lindeberg-Lévy)). *If  $E[X_1] = \mu$  and  $Var[X_1] = \sigma^2 > 0$  exist, then*

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow{\mathcal{L}} Z \sim N(0, 1).$$

**Theorem 3.2.3** (Multivariate CLT). *Let  $\{X_i\}$  be  $k$ -dim iid random vectors with mean  $\mu$  and covariance matrix  $\Sigma$ , then*

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{L}} Z \sim N_k(0, \Sigma).$$

**Theorem 3.2.4** (Strong Law of Large Numbers (SLLN, [6])). *If  $E[X_1] = \mu$  exists, then  $\bar{X}_n \xrightarrow{a.s.} \mu$ .*

The proofs of Theorems 3.2.1–3.2.3 using characteristic functions, are given in Section A.4.

Two proofs of Theorem 3.2.4 will be given in Section 3.3, one under a stronger condition that  $X_1$  has a finite fourth moment and the other assuming finite variance.

We conclude this section with a number of basic facts which follow from the definitions of various modes of convergence. We shall sketch the proofs of some of these, while the proofs of others are left as exercises. These facts will be used in the sequel without further explanation.

**Theorem 3.2.5.** *The following are true.*

- I. *Convergence in quadratic mean implies convergence in probability.*
- II. *If  $X_n \xrightarrow{P} X$  and  $g$  is continuous, then  $g(X_n) \xrightarrow{P} g(X)$ .*
- III. *If  $X_n \xrightarrow{\mathcal{L}} X$  and  $g$  is continuous (more generally, if  $P[X \in D_g] = 0$ , where  $D_g$  is the set of discontinuity points of  $g$ ), then  $g(X_n) \xrightarrow{\mathcal{L}} g(X)$ .*
- IV. (a)  *$X_n \xrightarrow{\mathcal{L}} X$  implies  $X_n = O_P(1)$ .*  
 (b)  *$X_n \xrightarrow{\mathcal{L}} 0$  implies  $X_n = o_P(1)$ .*
- V.  *$X_n \xrightarrow{\mathcal{L}} X$  and  $Y_n \xrightarrow{P} c$  implies  $X_n + Y_n \xrightarrow{\mathcal{L}} X + c$ ,  $X_n Y_n \xrightarrow{\mathcal{L}} cX$ , and if  $c \neq 0$ ,  $X_n/Y_n \xrightarrow{\mathcal{L}} X/c$ .*
- VI. (a) *If  $X_n = O_P(1)$  and  $Y_n = o_P(1)$ , then  $X_n Y_n = o_P(1)$  and  $X_n + Y_n = O_P(1)$ .*  
 (b) *If  $X_n = o_P(1)$  and  $Y_n = o_P(1)$ , then  $X_n Y_n = o_P(1)$  and  $X_n + Y_n = o_P(1)$ .*
- VII. *Slutsky's Theorem:  $X_n \xrightarrow{\mathcal{L}} X$  and  $Y_n = o_P(1)$  implies  $X_n + Y_n \xrightarrow{\mathcal{L}} X$ .*
- VIII. *Polya's Theorem: If  $F_n$  and  $F$  are cdf's,  $F$  is continuous, and  $F_n(x) \rightarrow F(x)$  for all  $x$ , then the convergence is uniform.*
- IX. *Borel-Cantelli Lemma: If  $\sum_{n=1}^{\infty} P[|X_n - a| > \varepsilon] < \infty$  for all  $\varepsilon > 0$ , then  $X_n \xrightarrow{a.s.} a$ . (The converse also holds under further conditions.)*
- X. (a) *If  $X_n \xrightarrow{a.s.} a$ , then  $X_n \xrightarrow{P} a$ .*  
 (b) *If  $X_n \xrightarrow{P} a$ , then there exists a sequence  $\{n_j\}$  such that the subsequence  $X_{n_j} \xrightarrow{a.s.} a$ .*

*Outlines of some proofs.*

- I. If  $E[(X_n - a)^2] \rightarrow 0$  as  $n \rightarrow \infty$ , then by Tchebyshev's inequality (to be proved in the next section), for any  $\varepsilon > 0$ ,

$$P[|X_n - a| > \varepsilon] \leq \frac{E[|X_n - a|^2]}{\varepsilon^2} \rightarrow 0, \text{ as } n \rightarrow \infty.$$

- II. (i) For  $\delta > 0$ , there exists  $M > 0$  such that  $P[|X| > M] \leq \delta$ .  
 (ii) For  $\delta > 0$ , there exists a positive integer  $n(\delta)$  such that  $P[|X_n - X| > \varepsilon] \leq \delta$  for all  $n \geq n(\delta)$ .

**(iii)** By (i) and (ii), for  $\delta > 0$ , there exist  $M > 0$  and  $n(\delta)$ , so that for all  $n \geq n(\delta)$ ,

$$P[|X| \leq M, |X_n - X| \leq \varepsilon, |X_n| \leq M + \varepsilon] \geq 1 - 2\delta,$$

**(iv)**  $g$  is continuous on  $\mathbb{R}$ , so  $g$  is uniformly continuous on  $[-M - \varepsilon, M + \varepsilon]$ .

Therefore, for  $\varepsilon_1 > 0$ , there exists  $\varepsilon_2 > 0$ , so that for  $x, x' \in [M - \varepsilon, M + \varepsilon]$ ,

$|x - x'| < \varepsilon_2$  implies  $|g(x) - g(x')| < \varepsilon_1$ .

**(v)** In (iii), choose  $\varepsilon = \varepsilon_2$  and choose  $\delta$  accordingly. Then

$$P[|g(X_n) - g(X)| < \varepsilon_1] \geq 1 - 2\delta \text{ for } n \geq n(\delta).$$

**III.** The proof is in [Section A.4](#).

**IV.** The proof of part (a) is in (iii) of II above. To prove part (b), note that  $X_n \xrightarrow{\mathcal{L}} 0$  means  $X_n \xrightarrow{\mathcal{L}} X$  where  $P[X = 0] = 1$ . Now taking  $A = (-\varepsilon, \varepsilon)$  in the definition of  $X_n \xrightarrow{\mathcal{L}} X$ , we have

$$P[|X_n| < \varepsilon] = P[X_n \in A] \rightarrow P[X \in A] = P[|X| < \varepsilon] = 1.$$

**V.** For any constant  $a$

$$P[X_n + Y_n \leq a] = P[X_n + Y_n \leq a, |Y_n - c| \leq \varepsilon] + P[X_n + Y_n \leq a, |Y_n - c| > \varepsilon].$$

Hence as  $n \rightarrow \infty$ ,

$$P[X_n + Y_n \leq a] \leq P[X_n \leq a - c + \varepsilon] + P[|Y_n - c| > \varepsilon] \rightarrow P[X \leq a - c + \varepsilon].$$

On the other hand

$$\begin{aligned} P[X_n + Y_n \leq a] &\geq P[X_n \leq a - c - \varepsilon, |Y_n - c| \leq \varepsilon] + 0 \\ &\geq P[X_n \leq a - c - \varepsilon] - P[|Y_n - c| > \varepsilon] \\ &\rightarrow P[X \leq a - c + \varepsilon]. \end{aligned}$$

Thus for any  $\varepsilon > 0$ ,

$$P[X \leq a - c - \varepsilon] \leq \lim_{n \rightarrow \infty} P[X_n + Y_n \leq a] \leq P[X \leq a - c + \varepsilon].$$

Hence

$$\lim_{n \rightarrow \infty} P[X_n + Y_n \leq a] = P[X \leq a - c] = P[X + c \leq a]$$

at all continuity points  $a$  of  $X + c$ . Hence  $X_n + Y_n \xrightarrow{\mathcal{L}} X + c$ .

**VII.** Slutsky's Theorem is a special case of (V) for  $c = 0$ .

**VIII.** To prove Polya's Theorem, choose  $x_1 < \dots < x_k$  such that  $F(x_1) < \varepsilon/2$ ,  $1 - F(x_k) < \varepsilon/2$ , and  $F(x_{i+1}) - F(x_i) < \varepsilon/2$ ,  $i = 1, \dots, k-1$ . Now find  $N$  such that  $|F_n(x_i) - F(x_i)| < \varepsilon/2$ ,  $i = 1, \dots, k$ , for  $n \geq N$ . The proof follows from these.

**IX.** Almost Sure Convergence ( $\xrightarrow{a.s.}$ ) takes place in the context of infinite sequences  $\{X_n, n = 1, 2, \dots\}$ . For this, the appropriate probability space is constructed by extending the probability spaces of  $(X_1, \dots, X_n)$ ,  $n = 1, 2, \dots$  in a suitable manner.

Now  $\lim_{n \rightarrow \infty} X_n$  is to be thought of in the space  $(\Omega, \mathcal{F}, P)$ , where  $\Omega$  is the space of all infinite sequences  $\omega = (x_1, x_2, \dots)$ ,  $\mathcal{F}$  is a suitable  $\sigma$ -field consisting of all sets  $C_n(B_n) = \{\omega: (x_1, x_2, \dots) : (x_1, \dots, x_n) \in B_n\}$ ,  $n = 1, 2, \dots$  and all Borel sets  $B_n$  in  $\mathbb{R}^n$  with  $P[C_n(B_n)] = P[(X_1, \dots, X_n) \in B_n]$ .

Almost sure convergence  $X_n \xrightarrow{a.s.} a$  has the following meaning in this space: for every  $\varepsilon > 0$  and for  $\omega \notin N$  where  $P[N] = 0$ , there exists  $n(\omega, \varepsilon)$  such that  $|X_n(\omega) - a| < \varepsilon$  for all  $n \geq n(\omega, \varepsilon)$ . If there is no such  $n(\omega, \varepsilon)$ , then the event  $A_n(\varepsilon) = \{\omega: |X_n(\omega) - a| \geq \varepsilon\}$  must occur infinity often (i.o.). Thus  $X_n \xrightarrow{a.s.} a$ , that is,  $P[\lim_{n \rightarrow \infty} X_n = a] = 1$  if  $P[A_n(\varepsilon) \text{ i.o.}] = 0$  for all  $\varepsilon > 0$ .

Now  $\{A_n(\varepsilon) \text{ i.o.}\}^c$  means that  $A_n(\varepsilon)^c$  must occur for all  $n \geq k$ , for some positive integer  $k$ , that is,

$$\begin{aligned} \{A_n(\varepsilon) \text{ i.o.}\}^c &= \bigcup_{k=1}^{\infty} \bigcap_{n=k}^{\infty} A_n^c(\varepsilon), \text{ so} \\ \{A_n(\varepsilon) \text{ i.o.}\} &= \left[ \bigcup_{k=1}^{\infty} \bigcap_{n=k}^{\infty} A_n^c(\varepsilon) \right]^c = \bigcap_{k=1}^{\infty} \bigcup_{n=k}^{\infty} A_n(\varepsilon) \end{aligned}$$

by De Morgan's rule (Chapter 1). Since  $\{\bigcup_{n=k}^{\infty} A_n(\varepsilon)\}$  is a decreasing sequence of sets,

$$P[|X_n - a| \geq \varepsilon \text{ i.o.}] = P[A_n(\varepsilon) \text{ i.o.}] = \lim_{k \rightarrow \infty} P\left[\bigcup_{n=k}^{\infty} A_n(\varepsilon)\right] \leq \lim_{k \rightarrow \infty} \sum_{n=k}^{\infty} P[A_n(\varepsilon)].$$

Hence  $\sum_{n=1}^{\infty} P[|X_n - a| \geq \varepsilon] = \sum_{n=1}^{\infty} P[A_n(\varepsilon)] < \infty$  implies  $\lim_{k \rightarrow \infty} \sum_{n=k}^{\infty} P[A_n(\varepsilon)]$ , so

$$P[|X_n - a| \geq \varepsilon \text{ i.o.}] = 0 \text{ for all } \varepsilon > 0,$$

that is,  $X_n \xrightarrow{a.s.} a$ . This proves the direct part of Borel-Cantelli Lemma.

**X(a).** Note that  $X_n \xrightarrow{a.s.} a$  iff  $\lim_{k \rightarrow \infty} P[\bigcup_{n=k}^{\infty} A_n(\varepsilon)] = 0$  for all  $\varepsilon > 0$ , which implies  $\lim_{k \rightarrow \infty} P[A_k(\varepsilon)] = \lim_{k \rightarrow \infty} P[|X_k - a| \geq \varepsilon] = 0$  for all  $\varepsilon > 0$ , that is,  $X_k \xrightarrow{P} a$ .

*Algebra of  $o_P$ ,  $O_P$ ,  $o$ , and  $O$ .* The properties (VIa, b) given above can be stated as:

$$\begin{aligned} O_P(1) + o_P(1) &= O_P(1), O_P(1)o_P(1) = o_P(1), \\ o_P(1) + o_P(1) &= o_P(1), \text{ and } o_P(1)o_P(1) = o_P(1). \end{aligned}$$

The following also hold

$$\begin{aligned} O_P(1) + o(1) &= O_P(1), O(1) + o_P(1) = O_P(1), O_P(1)o(1) = o_P(1), \\ O(1)o_P(1) &= o_P(1), o_P(1) + o(1) = o_P(1), \text{ and } o_P(1)o(1) = o_P(1). \end{aligned}$$

**Theorem 3.2.6** (Delta Method). Suppose that in  $\mathbb{R}^k$ ,  $\sqrt{n}(\mathbf{X}_n - \boldsymbol{\mu}) \xrightarrow{\mathcal{L}} \mathbf{Z} \sim N_k(\mathbf{0}, \boldsymbol{\Sigma})$  and  $g: \mathbb{R}^k \rightarrow \mathbb{R}$  has continuous first partial derivatives. Then

$$\sqrt{n}[g(\mathbf{X}_n) - g(\boldsymbol{\mu})] \xrightarrow{\mathcal{L}} W \sim N\left(0, [\nabla g(\boldsymbol{\mu})]^\top \boldsymbol{\Sigma} [\nabla g(\boldsymbol{\mu})]\right).$$

The proof of this theorem is left as an exercise.

In many instances of asymptotic normality, the variance of the asymptotic normal distribution is a function of its mean. This poses a problem in the construction of a large sample confidence interval for the mean, which is the parameter of interest. In such cases, it is convenient to make a transformation so that in the asymptotic distribution of the transform, the variance is a constant. These are called *variance-stabilizing transformations* of which some well-known examples are given below.

**Example 3.2.1.** If  $X$  is  $Poi(\mu)$ , then for large  $\mu$ ,  $(X - \mu)$  is asymptotically  $N(0, \mu)$ . Here  $(\sqrt{X} - \sqrt{\mu})$  is asymptotically  $N(0, 1/4)$ .

**Example 3.2.2.** If  $X_1, \dots, X_n$  are iid  $Bernoulli(p)$  and  $\hat{p}_n = n^{-1} \sum_{i=1}^{\infty} X_i$ , then  $\sqrt{n}(\hat{p}_n - p) \xrightarrow{\mathcal{L}} N(0, p(1-p))$ . Here  $\sqrt{n}(\arcsin \sqrt{\hat{p}_n} - \arcsin \sqrt{p}) \xrightarrow{\mathcal{L}} N(0, 1/4)$ .

**Example 3.2.3.** Let  $S_n^2 = (n-1)^{-1} \sum_{i=1}^{\infty} (X_i - \bar{X}_n)^2$ , where  $X_1, \dots, X_n$  are iid  $N(\mu, \sigma^2)$ . Then  $\sqrt{n}(S_n^2 - \sigma^2) \xrightarrow{\mathcal{L}} N(0, 2\sigma^4)$ . Here  $\sqrt{n}(\log S_n^2 - \log \sigma^2) \xrightarrow{\mathcal{L}} N(0, 2)$ .

**Example 3.2.4.** If  $(X_1, Y_1), \dots, (X_n, Y_n)$  are iid bivariate normal with  $\text{Corr}[X_i, Y_i] = \rho \in (-1, 1)$  and if  $r_n$  is the sample correlation coefficient, then  $\sqrt{n}(r_n - \rho) \xrightarrow{\mathcal{L}} N(0, (1 - \rho^2)^2)$ .

Here  $\sqrt{n}(\tanh^{-1} r_n - \tanh^{-1} \rho) \xrightarrow{\mathcal{L}} N(0, 1)$  (recall that  $\tanh^{-1} x = (1/2) \log[(1+x)/(1-x)]$ ).

Actually, using  $\sqrt{n-3}(\tanh^{-1} r_n - \tanh^{-1} \rho)$  as an asymptotically  $N(0, 1)$  rv for the purpose of constructing confidence intervals for  $\rho$  results in a better approximation.

To prove these results, find a transformation  $g$  in each case, such that  $\sqrt{n}(X_n - \mu) \xrightarrow{\mathcal{L}} N(0, \sigma^2(\mu))$  leads to  $\sqrt{n}[g(X_n) - g(\mu)] \xrightarrow{\mathcal{L}} N\left(0, \{g'(\mu)\}^2 \sigma^2(\mu)\right)$ , where  $\{g'(\mu)\}^2 \sigma^2(\mu)$  is a constant.

### 3.3 Probability Inequalities

In the previous section we have seen that the proofs of many “in probability” or “almost sure” convergence results require good upper bounds for tail probabilities of deviations of rv’s from their means, such  $P[|\bar{X}_n - \mu| \geq a]$  where  $\bar{X}_n$  is the mean of iid  $X_1, \dots, X_n$  with mean  $\mu$ . The following probability inequalities are useful for this purpose.

**Theorem 3.3.1** (Markov Inequality). If  $X$  is an rv with  $P[X \geq 0] = 1$ , then

$$P[X \geq a] \leq \frac{E[X]}{a}, \quad \text{for any } a > 0.$$

*Proof.* Note that for any  $a > 0$ ,  $I_{[a,\infty)}(x) \leq x/a$  and hence

$$P[X \geq a] = E[I_{[a,\infty)}(X)] \leq \frac{E[X]}{a}.$$

(See [Proposition 1.8.1](#)(iii and vii).)

□

**Tchebyshev's Inequality.** If  $X$  is an rv with mean  $\mu$  and variance  $\sigma^2$ , then

$$P[|X - \mu| \geq t] \leq \frac{\sigma^2}{t^2}, \quad \text{for any } t > 0.$$

*Proof.* Using Markov's inequality

$$P[|X - \mu| \geq t] = P[|X - \mu|^2 \geq t^2] \leq \frac{E[|X - \mu|^2]}{t^2} = \frac{\sigma^2}{t^2}.$$

□

*Generalization.* It follows in the same way, that if  $X$  is an rv with mean  $\mu$  and  $(2r)$ th central moment  $\mu_{2r} = E[|X - \mu|^{2r}]$ , then

$$P[|X - \mu| \geq t] \leq \frac{\mu_{2r}}{t^{2r}}, \quad \text{for any } t > 0.$$

□

*Applications.*

1. A simple proof of the WLLN assuming finite variance:

*WLLN.* If  $\bar{X}_n$  is the sample mean of iid rv's  $X_1, \dots, X_n$  with mean  $\mu$  and finite variance  $\sigma^2$ , then by Tchebyshev's inequality

$$\lim_{n \rightarrow \infty} P[|\bar{X}_n - \mu| \geq \varepsilon] \leq \lim_{n \rightarrow \infty} \frac{\text{Var}[\bar{X}_n]}{\varepsilon^2} = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\varepsilon^2} = 0.$$

2. A simple proof of the SLLN assuming finite fourth moment:

*SLLN.* If  $\bar{X}_n$  is the sample mean of iid rv's  $X_1, \dots, X_n$  with mean  $\mu$  and finite fourth central moment  $\tau^4$ , then by generalized Tchebyshev's inequality;

$$P[|\bar{X}_n - \mu| \geq \varepsilon] = P[|S_n| \geq n\varepsilon] \leq \frac{E[S_n^4]}{(n\varepsilon)^4},$$

where  $S_n = \sum_{i=1}^{\infty} (X_i - \mu) = \sum_{i=1}^{\infty} Y_i$ , where  $Y_1, \dots, Y_n$  are iid with mean 0 and finite fourth moment  $E[Y_i^4] = \tau^4$ . Therefore

$$\begin{aligned} E[S_n^4] &= E[(S_{n-1} + Y_n)^4] = E[S_{n-1}^4] + 6(n-1)\sigma^4 + \tau^4 = \dots \\ &= E[S_1^4] + 6\{1 + 2 + \dots + (n-1)\}\sigma^4 + (n-1)\tau^4 \\ &= 3n(n-1)\sigma^4 + n\tau^4. \end{aligned}$$

Thus

$$P[|\bar{X}_n - \mu| \geq \varepsilon] \leq \frac{3n(n-1)\sigma^4 + n\tau^4}{(n\varepsilon)^4} < \frac{3\sigma^4 + \tau^4}{n^2\varepsilon^4},$$

so that

$$\sum_{n=1}^{\infty} P[|\bar{X}_n - \mu| \geq \varepsilon] \leq \frac{3\sigma^4 + \tau^4}{\varepsilon^4} \sum_{n=1}^{\infty} \frac{1}{n^2} < \infty,$$

and now  $\bar{X}_n \xrightarrow{a.s.} \mu$  by the Borel-Cantelli Lemma.

This covers the case of many standard distributions where the SLLN holds. Another proof of the SLLN, assuming only finite variance, will be given later in this section, using a much stronger probability inequality.

We now obtain two more powerful inequalities which provide *exponential bounds* for tail probabilities. Of these, the first holds for bounded rv's and the second requires some moment conditions, providing a sharper bound. However, the first inequality due to Hoeffding has the advantage of simplicity and is very useful in many situations.

**Theorem 3.3.2** (Hoeffding's Inequality [7]). *If  $X_1, \dots, X_n$  are independent with  $P[0 \leq X_i \leq 1] = 1$  for all  $i$ ,  $S_n = X_1 + \dots + X_n$ ,  $\bar{X}_n = S_n/n$ , and  $\mu = E[\bar{X}_n] = E[S_n]/n$ , then for  $0 \leq t \leq 1 - \mu$ ,  $P[\bar{X}_n - \mu \geq t] \leq e^{-2nt^2}$ .*

*Proof.* Note that

$$\begin{aligned} P[\bar{X}_n - \mu \geq t] &= P[S_n - n\mu - nt \geq 0] = E[I_{[0,\infty)}(S_n - n\mu - nt)] \\ &\leq E[e^{h(S_n - n\mu - nt)}], \text{ for } h > 0 \\ &= e^{-nh(\mu+t)} \prod_{i=1}^n E[e^{hX_i}] \leq e^{-nh(\mu+t)} \prod_{i=1}^n E[(1 - X_i)e^0 + X_i e^h] \\ &= e^{-nh(\mu+t)} \prod_{i=1}^n [1 - \mu_i + \mu_i e^h] \leq e^{-nh(\mu+t)} \left[ n^{-1} \sum_{i=1}^n (1 - \mu_i + \mu_i e^h) \right]^n \\ &= e^{-nh(\mu+t)} (1 - \mu + \mu e^h)^n \\ &= e^{nL(h)}, \text{ with } L(h) = \log(1 - \mu + \mu e^h) - h(\mu + t), \end{aligned}$$

using the convexity of  $e^{hx}$  and because geometric mean is no larger than the arithmetic mean. Next, we want to choose  $h > 0$  so that  $L(h) < 0$ . This is possible because  $L'(0) = -t$  and  $L''(h) \leq 1/4$  for all  $h$ , as can be easily verified, so that

$$\begin{aligned} L(h) &= L(0) + hL'(0) + \frac{1}{2}h^2L''(h^*) = 0 - ht + \frac{1}{2}h^2L''(h^*) \\ &\leq -ht + \frac{h^2}{8} = -2t^2 \text{ for } h = 4t. \end{aligned}$$

Hence we get the inequality by taking  $h = 4t$ . □

**Theorem 3.3.3** (Bernstein's Inequality). *If  $X_1, \dots, X_n$  are independent with  $E[X_i] = 0$ ,  $E[X_i^2] = b_i$ ,  $B_n = b_1 + \dots + b_n$ , and for  $r > 2$ ,*

$$E[|X_i|^r] \leq \frac{1}{2} r! b_i c^{r-2}, \quad i = 1, \dots, n,$$

where  $c$  is a constant, then

$$P[|X_1 + \dots + X_n| > t] \leq 2 \exp\left[-\frac{t^2}{2B_n + 2ct}\right].$$

*Proof.* For an outline of the proof, see Uspensky [8, p. 204–5].  $\square$

**Corollary.** *If  $P[|X_i| \leq M] = 1$ , then for  $r > 2$ ,*

$$\begin{aligned} E[|X_i|^r] &= E[X_i^2 |X_i|^{r-2}] \leq b_i M^{r-2} \\ &< \frac{b_i M^{r-2} r!}{2 \cdot 3^{r-2}} = \frac{1}{2} r! b_i \left(\frac{M}{3}\right)^{r-2}, \end{aligned}$$

because  $r!/(2 \cdot 3^{r-2}) > 1$ . In Bernstein's inequality we can now take  $c = M/3$ , resulting in

$$P[|X_1 + \dots + X_n| > t] \leq 2 \exp\left[-\frac{t^2}{2B_n + 2Mt/3}\right].$$

In particular, if  $X_1, \dots, X_n$  are iid Bernoulli( $p$ ), then  $M = \max(p, 1-p)$ , and we have

$$P[|\bar{X}_n - p| > t] \leq 2 \exp\left[-\frac{1}{2} \cdot \frac{nt^2}{p(1-p) + \frac{t}{3} \max(p, 1-p)}\right].$$

Finally we obtain two bounds for tail probabilities of the maximum of cumulative sums of independent rv's. For this reason, they are called maximal inequalities.

In what follows,  $X_1, \dots, X_n$  is a sequence of independent rv's with  $E[X_i] = 0$ ,  $E[X_i^2] = \sigma_i^2$ ,  $S_k = X_1 + \dots + X_k$ , and  $\varepsilon > 0$ .

**Theorem 3.3.4** (Kolmogorov's Inequality).

$$P\left[\max_{1 \leq k \leq n} |S_k| \geq \varepsilon\right] \leq \frac{1}{\varepsilon^2} \sum_{1 \leq k \leq n} \sigma_k^2.$$

*Proof.* In the space of infinite sequences  $\omega = (x_1, x_2, \dots)$ , let  $A_1 = \{\omega: |X_1| \geq \varepsilon\}$  and  $A_r = \{\omega: |S_k| < \varepsilon, k = 1, \dots, r-1, |S_r| \geq \varepsilon\}$ , for  $r \geq 2$ , that is,  $A_r$  is the set of all  $\omega$  for which the cumulative sums  $S_1, S_2, \dots$  go beyond  $\pm\varepsilon$  for the first time. Then  $A_1, A_2, \dots$  are mutually exclusive and

$$\left\{\omega: \max_{1 \leq k \leq n} |S_k| \geq \varepsilon\right\} = \bigcup_{k=1}^n A_k \quad \text{so } P\left[\max_{1 \leq k \leq n} |S_k| \geq \varepsilon\right] = \sum_{k=1}^n P[A_k].$$

Now

$$\begin{aligned}
\sum_{k=1}^n \sigma_k^2 &= E[(X_1 + \cdots + X_n)^2] \geq \sum_{k=1}^n P[A_k] E[(X_1 + \cdots + X_n)^2 | A_k] \\
&= \sum_{k=1}^n P[A_k] \left\{ E[(X_1 + \cdots + X_k)^2 | A_k] + E[(X_{k+1} + \cdots + X_n)^2 | A_k] \right. \\
&\quad \left. + 2 \sum_{i=1}^k \sum_{j=k+1}^n E[X_i X_j | A_k] \right\} \\
&\geq \sum_{k=1}^n P[A_k] \{ \varepsilon^2 + 0 + 0 \} = \varepsilon^2 \sum_{k=1}^n P[A_k] = \varepsilon^2 P \left[ \max_{1 \leq k \leq n} |S_k| \geq \varepsilon \right].
\end{aligned}$$

Hence the inequality follows. In the above argument,

- (i)  $E[(X_1 + \cdots + X_k)^2 | A_k] \geq \varepsilon^2$  because  $A_k$  is defined that way, and
- (ii) for  $1 \leq i \leq k, k+1 \leq j \leq n$ ,

$$\begin{aligned}
E[X_i X_j | A_k] &= E[E(X_i X_j | X_1, \dots, X_k) | A_k] \\
&= E[X_i E(X_j | X_1, \dots, X_k) | A_k] \\
&= E[X_i E(X_j) | A_k] = 0,
\end{aligned}$$

since the  $X_i$ 's are independent and  $E(X_j) = 0$ .  $\square$

*Remark 3.3.1.* In this proof, independence is used only in (ii) which holds more generally, so long as  $E[X_j | X_1, \dots, X_k] = 0$  for  $j \geq k+1$ . Let  $\mathcal{F}_k$  be the  $\sigma$ -field generated by  $X_1, \dots, X_k$ . Then  $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$  and  $E[S_j | \mathcal{F}_k] = S_k$  for  $j \geq k+1$ . This property of  $(S_k, \mathcal{F}_k)_{1 \leq k < \infty}$  is called the *martingale property*, which is enough for this theorem (see [9, p. 105]).

**Theorem 3.3.5** (Hájek-Rényi Inequality). *Let  $c_1, c_2, \dots$  be a nonincreasing sequence of positive numbers. Then for any two positive integers  $m < n$ ,*

$$P \left[ \max_{m \leq k \leq n} c_k | S_k | \geq \varepsilon \right] \leq \frac{1}{\varepsilon^2} \left[ c_m^2 \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^n c_k^2 \sigma_k^2 \right].$$

(For  $m = 1$  and  $c_1 = \dots = c_n$ , this reduces to the Kolmogorov inequality.)

*Proof.* Generalizing the definition of  $\{A_k\}$  in the proof of the Kolmogorov inequality, let

$$A_m = \{\omega: c_m |S_m| \geq \varepsilon\} \text{ and } A_r = \{\omega: c_k |S_k| < \varepsilon, m \leq k \leq r-1, c_r |S_r| \geq \varepsilon\}$$

for  $r \geq m+1$ , and as a generalization of  $(X_1 + \cdots + X_n)^2$ , let

$$Z = \sum_{k=m}^{n-1} (c_k^2 - c_{k+1}^2) S_k^2 + c_n^2.$$

Then

$$\begin{aligned}
 E[Z] &= \sum_{k=m}^{n-1} (c_k^2 - c_{k+1}^2) (\sigma_1^2 + \cdots + \sigma_k^2) + c_n^2 (\sigma_1^2 + \cdots + \sigma_n^2) \\
 &= \sum_{k=m}^{n-1} c_k^2 \sum_{j=1}^k \sigma_j^2 - \sum_{k=m+1}^n c_k^2 \sum_{j=1}^{k-1} \sigma_j^2 + c_n^2 \sum_{j=1}^n \sigma_j^2 \\
 &= \left\{ c_m^2 \sum_{j=1}^m \sigma_j^2 + \sum_{k=m+1}^{n-1} c_k^2 \sum_{j=1}^k \sigma_j^2 \right\} - \left\{ \sum_{k=m+1}^{n-1} c_k^2 \sum_{j=1}^{k-1} \sigma_j^2 + c_n^2 \sum_{j=1}^{n-1} \sigma_j^2 \right\} \\
 &\quad + c_n^2 \sum_{j=1}^n \sigma_j^2 \\
 &= c_m^2 \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^n c_k^2 \sigma_k^2, \text{ and}
 \end{aligned}$$

$$P\left[\max_{m \leq k \leq n} c_k |S_k| \geq \varepsilon\right] = \sum_{k=m}^n P[A_k].$$

We now proceed as in the previous proof:

$$\begin{aligned}
 c_m^2 \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^n c_k^2 \sigma_k^2 &= E[Z] \geq \sum_{k=m}^n P[A_k] E[Z|A_k] \\
 &= \sum_{k=m}^n P[A_k] E\left[\sum_{j=m}^{n-1} (c_j^2 - c_{j+1}^2) S_j^2 + c_n^2 S_n^2 | A_k\right] \\
 &\geq \sum_{k=m}^n P[A_k] E\left[\sum_{j=k}^{n-1} (c_j^2 - c_{j+1}^2) \left\{ S_k^2 + (S_j - S_k)^2 + 2 \sum_{r=1}^k \sum_{s=k+1}^j X_r X_s \right\} \right. \\
 &\quad \left. + c_n^2 \left\{ S_k^2 + (S_n - S_k)^2 + 2 \sum_{r=1}^k \sum_{s=k+1}^n X_r X_s \right\} | A_k \right] \\
 &\geq \sum_{k=m}^n P[A_k] \left[ \sum_{j=k}^{n-1} (c_j^2 - c_{j+1}^2) \frac{\varepsilon^2}{c_k^2} + \frac{c_n^2 \varepsilon^2}{c_k^2} \right] \\
 &= \varepsilon^2 \sum_{k=m}^n P[A_k] = \varepsilon^2 P\left[\max_{m \leq k \leq n} c_k |S_k| \geq \varepsilon\right].
 \end{aligned}$$

□

*Application of Hájek-Rényi Inequality.* A simple proof of the SLLN assuming finite variance is given here. Using the notations used in the above proof we see that for  $c_k = 1/k$ ,  $c_k |S_k| = |\bar{X}_k|$ . Hence by the H-R inequality,

$$\begin{aligned}
P\left[\max_{k \geq m} |\bar{X}_k| = \max_{k \geq m} c_k |S_k| \geq \varepsilon\right] &= \lim_{n \rightarrow \infty} P\left[\max_{m \leq k \leq n} c_k |S_k| \geq \varepsilon\right] \\
&\leq \lim_{n \rightarrow \infty} \frac{1}{\varepsilon^2} \left[ c_m^2 \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^n c_k^2 \sigma_k^2 \right] \\
&= \frac{1}{\varepsilon^2} \left[ c_m^2 \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^{\infty} c_k^2 \sigma_k^2 \right] \\
&= \frac{1}{\varepsilon^2} \left[ m^{-2} \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^{\infty} \frac{\sigma_k^2}{k^2} \right].
\end{aligned}$$

Now suppose that  $E[X_k^2] = \sigma^2$  for all  $k$ , that is,  $X_1, X_2, \dots$  have finite variance  $\sigma^2$ . Then

$$\lim_{m \rightarrow \infty} \left[ m^{-2} \sum_{k=1}^m \sigma_k^2 + \sum_{k=m+1}^{\infty} \frac{\sigma_k^2}{k^2} \right] = \lim_{m \rightarrow \infty} \sigma^2 \left[ \frac{m}{m^2} + \sum_{k=m+1}^{\infty} \frac{1}{k^2} \right] = 0.$$

Hence  $\lim_{m \rightarrow \infty} P[\max_{k \geq m} |\bar{X}_k| \geq \varepsilon] = 0$ .

But the events

$$\left\{ \omega: \max_{k \geq m} |\bar{X}_k| \geq \varepsilon \right\} = \bigcup_{k=m}^{\infty} \{\omega: |\bar{X}_k| \geq \varepsilon\} = \bigcup_{k=m}^{\infty} B_k(\varepsilon)$$

are nonincreasing as  $m$  increases. Therefore

$$\begin{aligned}
0 &= \lim_{m \rightarrow \infty} P\left[\max_{k \geq m} |\bar{X}_k| \geq \varepsilon\right] = P\left[\bigcap_{m=1}^{\infty} \bigcup_{k=m}^{\infty} B_k(\varepsilon)\right] \\
&= P[B_k(\varepsilon) \text{ i.o.}] = P[|\bar{X}_k| \geq \varepsilon \text{ i.o.}] = P[\bar{X}_k \not\rightarrow 0].
\end{aligned}$$

Thus  $\bar{X}_k \xrightarrow{a.s.} 0$  if  $\{X_k\}$  is a sequence of independent rv's with mean 0 and finite variance  $\sigma^2$ .

### 3.4 Asymptotic Normality: The Central Limit Theorem and Its Generalizations

The basic Lindeberg-Lévy CLT has been stated in [Theorem 3.2.2](#) and its proof, using characteristic functions is given in [Section A.4](#). The Lindeberg-Lévy CLT is restricted to sample means of iid rv's with finite variance. However, in statistical inference based on large samples, we often have to deal with sample means of *independent but nonidentically distributed* rv's and in some situations, even the distribution of the  $n$  independent rv's changes from one  $n$  to another, as well as the sample size  $n$ . In the following, we state two generalizations of [Theorem 3.2.2](#) without proof which address these generalities.

**Theorem 3.4.1** (Lindeberg-Feller Theorem). *Let  $X_1, X_2, \dots$  be independent rv's with  $E[X_i] = \mu_i$ ,  $\text{Var}[X_i] = \sigma_i^2 < \infty$ , and suppose  $B_n^2 = \sigma_1^2 + \dots + \sigma_n^2$  satisfies  $\lim_{n \rightarrow \infty} B_n = \infty$  and  $\lim_{n \rightarrow \infty} (\sigma_n^2/B_n^2) = 0$ . Then*

$$B_n^{-1} \sum_{i=1}^n (X_i - \mu_i) \xrightarrow{\mathcal{L}} Z \sim N(0, 1) \text{ iff}$$

$$\lim_{n \rightarrow \infty} B_n^{-2} \sum_{i=1}^n E[I_{(\varepsilon B_n, \infty)}(|X_i - \mu_i|)(X_i - \mu_i)^2] = 0 \quad \text{for all } \varepsilon > 0.$$

*Triangular Arrays.* Here we think of a double array of rv's in which each row consists of independent rv's

$$\{(X_{11}, \dots, X_{1k_1}), (X_{21}, \dots, X_{2k_2}), \dots, (X_{n1}, \dots, X_{nk_n}), \dots\}$$

with  $E[X_{ni}] = \mu_{ni}$ ,  $\text{Var}(X_{ni}) = \sigma_{ni}^2 < \infty$ , and  $B_n^2 = \sigma_{n1}^2 + \dots + \sigma_{nk_n}^2$ .

**Theorem 3.4.2** (Lindeberg-Liapounov Theorem). *The Lindeberg condition*

$$\lim_{n \rightarrow \infty} B_n^{-2} \sum_{i=1}^{k_n} E[I_{(\varepsilon B_n, \infty)}(|X_{ni} - \mu_{ni}|)(X_{ni} - \mu_{ni})^2] = 0 \quad \text{for all } \varepsilon > 0$$

*implies  $B_n^{-1} \sum_{1 \leq i \leq k_n} (X_{ni} - \mu_{ni}) \xrightarrow{\mathcal{L}} Z \sim N(0, 1)$ . The above convergence in law is also implied by the Liapounov condition*

$$\lim_{n \rightarrow \infty} \frac{\rho_n^3}{B_n^3} = 0, \text{ where } \rho_n^3 = \sum_{i=1}^{k_n} E[|X_{ni} - \mu_{ni}|^3] < \infty \text{ and } B_n^3 = \left( \sum_{i=1}^{k_n} \sigma_{ni}^2 \right)^{3/2}.$$

Finally, we state another generalization of [Theorem 3.2.2](#) for sequences of rv's with limited dependence.

**Definition 3.4.1.** A sequence of rv's  $\{X_1, X_2, \dots\}$  is said to be  $m$ -dependent if  $(X_1, \dots, X_r)$  is independent of  $(X_s, X_{s+1}, \dots)$  whenever  $s - r > m$ .

Let  $A_i = 2 \sum_{j=0}^{m-1} \text{Cov}[X_{i+j}, X_{i+m}] + \text{Var}[X_{i+m}]$ .

**Theorem 3.4.3** (Hoeffding and Robbins [10]). *If for an  $m$  dependent sequence  $\{X_i\}$  with  $E[X_i] = 0$ ,  $\text{Var}[|X_i|^3] \leq K < \infty$ , and  $\lim_{r \rightarrow \infty} r^{-1} \sum_{j=1}^r A_{i+j} = A > 0$  exists uniformly for all  $i$ , then  $n^{-1/2} \sum_{i=1}^n X_i \xrightarrow{\mathcal{L}} N(0, A)$ .*

*In particular, if  $\{X_i\}$  is stationary  $m$ -dependent sequence with  $E[X_i] = \mu$ ,  $E[|X_i|^3] < \infty$  and  $A = \text{Var}[X_1] + 2 \sum_{i=2}^{m+1} \text{Cov}[X_1, X_i] > 0$ , then  $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{L}} N(0, A)$ .*

## Exercises

For all the problems below, when there are observations  $X_1, \dots, X_n$ , it is understood that

$$\bar{X}_n = n^{-1} \sum_{i=1}^n X_i, \quad s_n^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

$$X_{n:1} = \min(X_1, \dots, X_n), \quad X_{n:n} = \max(X_1, \dots, X_n).$$

- 3.1. Prove VI(a,b) of [Theorem 3.2.5](#).
- 3.2. Prove X(b) of [Theorem 3.2.5](#).
- 3.3. Prove [Theorem 3.2.6](#), the Delta Method.
- 3.4. Prove the results stated in the examples on applications of the Delta Method.
- 3.5. Let  $X_1, \dots, X_n$  be iid with  $E[X_i] = \mu \neq 0$ ,  $\text{Var}[X_i] = \sigma^2$  and  $E[X_i^4] < \infty$ . Find the asymptotic distribution of  $\sqrt{n}(s_n/\bar{X}_n - \sigma/\mu)$  as  $n \rightarrow \infty$ .
- 3.6. Let  $B_{m,n}$  denote a  $\text{Beta}(m, n)$  rv. Show that if  $m, n \rightarrow \infty$  in such a way that  $m/(m+n) \rightarrow \alpha \in (0, 1)$ , then  $\sqrt{m+n}(B_{m,n} - m/(m+n))/\sqrt{\alpha(1-\alpha)}$  converges in distribution to  $Z \sim N(0, 1)$ . [First show that if  $X_1, \dots, X_m, Y_1, \dots, Y_n$  are iid  $\text{Exp}(1)$ , then  $R_{m,n}(1 + R_{m,n})^{-1} \xrightarrow{D} B_{m,n}$ , where  $R_{m,n} = (m\bar{X}_m)/(n\bar{Y}_n)$ .]
- 3.7. Let  $\{X_i\}$  be a sequence of rv's with  $E[X_i] = 0$ ,  $\text{Var}[X_i] = 1$  and  $\text{Cov}[X_i, X_j] = 0$  for  $|i - j| \geq k$ , where  $k$  is a fixed positive integer. Show that  $\bar{X}_n \xrightarrow{P} 0$ .
- 3.8. Let  $\{X_n\}$  be a sequence of rv's with  $E[X_n] = \mu_n$  and  $\text{Var}[X_n] = \sigma_n^2$ . Show that if  $\mu_n \rightarrow 0$  and  $\sigma_n^2 \rightarrow 0$  as  $n \rightarrow \infty$ , then  $X_n \xrightarrow{P} 0$ .
- 3.9. Let  $\{X_n\}$  be a sequence of rv's such that  $X_n$  is distributed as  $Z$  with probability  $p_n$  and as  $\sigma_n Z$  with probability  $1 - p_n$  where  $Z \sim N(0, 1)$ ,  $p_n \rightarrow p \in [0, 1]$  and  $\sigma_n \rightarrow \infty$ .
  - (a) Show that  $X_n \xrightarrow{\mathcal{L}} X$  if and only if  $p = 1$ . What is the distribution of  $X$ ?
  - (b)  $\lim_{n \rightarrow \infty} \text{Var}[X_n]$  is not necessarily the same as  $\text{Var}[X]$ . Find the possible values of  $\lim_{n \rightarrow \infty} \text{Var}[X_n]$ .
- 3.10. Let  $\{(x_{ni}, Y_{ni}), i = 1, \dots, n\}_{n=1}^\infty$  be a triangular array following a simple linear regression model  $Y_{ni} = \alpha + \beta x_{ni} + Z_{ni}$ , where  $Z_{n1}, \dots, Z_{nn}$  are iid with mean 0 and variance  $\sigma^2$ . Assume that the common distribution of  $Z_{ni}$ 's is the same for all  $n$ . Let

$$\hat{\beta}_n = \frac{\sum_{i=1}^n (x_{ni} - \bar{x}_n) Y_{ni}}{\sum_{i=1}^n (x_{ni} - \bar{x}_n)^2}, \quad \bar{x}_n = n^{-1} \sum_{i=1}^n x_{ni}$$

denote the least squares estimator of  $\beta$ . Use Lindeberg's condition to establish the asymptotic normality of  $\sqrt{n}(\hat{\beta}_n - \beta)$ , making appropriate assumptions on  $\{x_{ni}\}$ .

# Basic Concepts of Statistical Inference

## 4.1 Population and Random Samples

The term “Statistics” is commonly used as a synonym for data, but “Statistical Inference” is the science of analyzing data to probe into where the data came from. The data is a “sample” and where the data came from is the “population” having some unknown characteristics in which we are interested. The sample has to be a “random sample” for the sake of objectivity and thus randomness brings probability into the picture.

A population can be thought of as a concrete, finite collection, such as individuals in a city, from which a random sample is drawn and some characteristic of each of these individuals, such as opinion on a certain issue, or income, etc., is recorded. This is the data, from which a summary measure of the characteristic in the entire population must be inferred. This is the framework of Survey Sampling.

Here, on the other hand, we think of a population consisting of observations on the outcomes of infinite repetitions of an experiment, such as survival times of cancer patients on a certain drug, or the number of defective items in lots of  $N$  items coming out of a production line, or the number of accidents during certain hours on a particular stretch of a freeway. Clearly, these observations vary from one experiment to another and they vary in a random manner. The collection of possible outcomes of such an experiment is a probability space  $(S, \mathcal{A}, P)$  and the observation on an outcome,  $X = X(s)$  is a random variable or more generally, a random vector. From  $n$  independent repetitions of such an experiment, we observe independent random variables  $X_1, \dots, X_n$ . This is the data.

The probability distribution of  $X$  should be denoted by  $P^X$  in strict notation (to distinguish it from the set function  $P: \mathcal{A} \rightarrow [0, 1]$  in the basic probability space), defined as  $P^X[(-\infty, a]] = P[\{s: X(s) \leq a\}]$  and more generally by  $P^X[B] = P[\{s: X(s) \in B\}]$  for all Borel sets in  $\mathbb{R}$ . However, to keep the notation simple, we shall use  $P$  instead of  $P^X$  to denote the probability distribution of  $X$  when there is no chance of confusion. These notations will extend to the case of multidimensional rv's in an obvious manner. We can now say that our observations  $X_1, X_2, \dots$  are independent rv's with common probability distribution  $P$  which is an unknown element of a family of distributions  $\mathbb{P}$ .

## 4.2 Parametric and Nonparametric Models

Instead of leaving the family  $\mathbb{P}$  wide open, we shall assume that  $\mathbb{P}$  is either a *parametric family*

$$\mathbb{P} = \{P_\theta, \theta \in \Theta\},$$

where  $\theta = (\theta_1, \dots, \theta_k)$  is a  $k$ -dim parameter vector belonging to the parameter space  $\Theta$  and  $P_\theta$  is known for any given  $\theta \in \Theta$ , or a *nonparametric family*  $\mathbb{P}$  whose members cannot be identified by a finite number of parameters.

Examples of parametric family are

- P(i)** *Bernoulli* ( $\theta$ ),  $0 < \theta < 1$ ,
- P(ii)** *Poisson* ( $\lambda$ ),  $\lambda > 0$ ,
- P(iii)** *Normal* ( $\mu, \sigma^2$ ),  $-\infty < \mu < \infty$  and  $\sigma^2 > 0$ ,
- P(iv)** *Gamma* ( $\alpha, \beta$ ),  $\alpha > 0$  and  $\beta > 0$ ,
- P(v)** Linear Regression Model ( $\alpha, \beta, \sigma^2$ ), etc.

The families P(i)–P(iv) have been discussed in [Chapter 2](#) and P(v) has been introduced in [Chapter 1, Section 1.10](#). In P(v), the distributions of  $X$  and  $\varepsilon$  are not fully specified, but we still can do a lot within the model as it is and much more conditionally given  $X = x$  if  $\varepsilon \sim N(0, \sigma^2)$ .

Examples of nonparametric family are

- NP(i)** All probability distributions on the real line.
- NP(ii)** All probability distributions on the real line with pdf's satisfying some smoothness conditions.
- NP(iii)** The family described in NP(ii) with the further restriction that the pdf's are symmetric about some  $\theta \in \mathbb{R}$ .
- NP(iv)** All probability distribution of  $X_1$  and  $X_2$  on the real line which are independent and

$$f_{X_2}(x) = f_{X_1}(x - \theta) \quad \text{for all } x \text{ and for some } \theta \in \mathbb{R}.$$

- NP(v)** All probability distributions of  $(X_1, \dots, X_k, Y)$  on  $\mathbb{R}^{k+1}$  with  $f_{X_1, \dots, X_k}$  and  $m(x_1, \dots, x_k) = E[Y|X = \mathbf{x}]$ , satisfying some smoothness conditions.

## 4.3 Problems of Statistical Inference

Based on observed data consisting of random samples  $(X_1, \dots, X_n)$  or  $(X_{11}, \dots, X_{1n}; X_{21}, \dots, X_{2n})$  or  $((X_1, Y_1), \dots, (X_n, Y_n))$ , etc., from an unknown probability distribution  $P \in \mathbb{P}$  in a parametric or nonparametric model, we have to make inference about some unknown features of  $P$ . The three main types of inference that the statistical science has been concerned with from its very inception are *Point Estimation*, *Hypothesis Testing*, and

*Confidence Sets.* We briefly describe these three types here by means of some examples before taking them up in more details in subsequent chapters.

## Point Estimation

Here we want to make a guess about a function  $g(\theta)$  of the unknown  $\theta \in \Theta$  in a parametric model, or construct a pdf or a regression function with prescribed properties as our guess for an unknown pdf  $f$  or an unknown regression function  $m$  in a nonparametric model. In a parametric model, we may want to estimate the mean  $\lambda$  of a Poisson distribution or the slope  $\beta$  of a linear regression function or  $P[X \leq a]$  of  $X$  distributed as  $N(\mu, \sigma^2)$ .

## Hypothesis Testing

In a parametric model, this involves deciding whether  $\theta \in \Theta_0$  or  $\theta \in \Theta_1$  where  $\Theta_0$  and  $\Theta_1$  are disjoint sets in the parameter space  $\Theta$ , such as deciding in the context of  $N(\mu, \sigma^2)$  with  $\sigma^2$  known or unknown, whether  $\mu \leq 0$  or  $\mu > 0$ . We call  $H_0: \theta \in \Theta_0$  the *null hypothesis* which we are inclined to believe unless the data provides *significant* evidence in favor of  $H_1: \theta \in \Theta_1$  which we call the *alternative hypothesis*. So there is an asymmetry in this problem due to our attitude toward  $H_0$  and  $H_1$ . For this reason, this is a problem of testing  $H_0$  against  $H_1$ . An example in a nonparametric model is to decide whether or not  $\theta = 0$  (ie, testing  $H_0: \theta = 0$  against  $H_1: \theta \neq 0$ ) in the nonparametric family NP(iv).

## Confidence Sets

Unlike point estimation, here we want to *construct a set* in the parameter space to which we guess the unknown  $\theta$  to belong. In a parametric model, this may be an interval which includes the mean  $\mu$  of a normal distribution  $N(\mu, \sigma^2)$  with  $\sigma^2$  known or unknown, while in a nonparametric model this may be a band in which we guess that an unknown cdf  $F$  belongs.

These three problems will be treated in [Chapters 5 and 6](#) within parametric families and in [Chapter 8](#) within nonparametric families.

## 4.4 Statistical Decision Functions

Problems of statistical inference in a parametric model such as Point Estimation, Hypothesis Testing, and many others can be fitted in a general framework of decision making based on a random sample  $X$  from  $(\mathfrak{X}, \mathbb{B}, P_\theta)$  with unknown  $\theta \in \Theta$ . The decision consists of taking an action  $a \in A$ , where  $A$  is the set of available actions which we call the *action space*. This action has to be a function of the observed value  $x$  of  $X$ . We call this function a *Decision Function* or *Decision Rule*:

$$d: \mathfrak{X} \rightarrow A \text{ with } d(x) = a \in A \text{ when } X = x.$$

Now there are good decisions and bad decisions for any given  $\theta \in \Theta$ , and for making a bad decision one has to pay a penalty, or take a loss, which is the consequence of taking the action  $a = d(x)$  when the observed value  $x$  of  $X$  was generated from  $P_\theta$ . This leads to the definition of a *Loss Function*:  $L: \Theta \times A \rightarrow \mathbb{R}$  with  $L(\theta, a) :=$  loss due to action  $a$  when  $\theta$  is the true value of the parameter. Then for  $X = x$ ,  $L(\theta, d(x)) =$  loss due to action  $d(x)$  when  $X = x$  is generated by  $P(\theta)$ .

Since the data is the rv  $X$ , we should be thinking of the loss as an rv  $L(\theta, d(X))$  where  $X \sim P_\theta$ . Note that the  $\theta$  in  $P_\theta$  and the  $\theta$  in  $L(\theta, d(X))$  is the same unknown element of  $\Theta$  (ie, the true value of the parameter).

A statistical decision problem is thus described by the triple  $(\{P_\theta: \theta \in \Theta\}, A, L)$ , and the overall performance of a decision function  $d$  is measured by its risk

$$\begin{aligned} R(\theta, d) &= E_\theta[L(\theta, d(X))] \\ &= \begin{cases} \int_{\mathfrak{X}} L(\theta, d(x))f(x, \theta) dx & \text{in the continuous case,} \\ \sum_{x \in \mathfrak{X}} L(\theta, d(x))f(x, \theta) & \text{in the discrete case,} \end{cases} \end{aligned}$$

assuming that the  $\int$  or the  $\sum$  defining  $R(\theta, d)$  exists. This holds if  $\mathfrak{X}$  is finite and more generally, if  $L(\cdot, \cdot)$  is bounded below, although the  $\int$  or the  $\sum$  may be  $+\infty$ .

**Definition 4.4.1.** The function  $R(\cdot, d)$  is called the risk function of the decision rule  $d$ .

Ideally, we should be using a decision rule  $d^*$  for which  $R(\theta, d^*) \leq R(\theta, d)$ , for all  $\theta \in \Theta$  and for all decision rules  $d$ .

Unfortunately, such a decision rule does not exist (except in some trivial cases). This is illustrated by the following examples. Before introducing these examples, we enlarge the class of decision rules by allowing randomization over  $a \in A$  for each  $x$ .

**Definition 4.4.2.** A *behavioral decision rule*  $\delta$  consists of a probability distribution  $\delta(\cdot|x): \mathcal{C} \rightarrow [0, 1]$  for each  $x \in \mathfrak{X}$ , where  $\mathcal{C}$  is a  $\sigma$ -field of subsets of the action space  $A$ .

If  $A$  is countable, then  $\mathcal{C}$  could be the class of all subsets of  $A$ . For  $A = \mathbb{R}^k$  or some subsets of  $\mathbb{R}^k$ ,  $\mathcal{C}$  could be the Borel sets and  $\delta(\cdot|x)$  would have a density. The risk of a behavioral decision rule  $\delta(\cdot|x)$  is given by:

$$R(\theta, \delta) = \int_{\mathfrak{X}} \int_A L(\theta, a) d\delta(a|x) f(x, \theta) dx,$$

where  $\int_A L(\theta, a) d\delta(a|x) = \sum_{a \in A} L(\theta, a)\delta(a|x)$  if  $A$  is countable and

$$\int_A L(\theta, a) d\delta(a|x) = \int_A L(\theta, a)p(a|x) da$$

when  $A = \mathbb{R}^k$  is some subset of  $\mathbb{R}^k$ ,  $p(\cdot|x)$  being the pdf of  $\delta(\cdot|x)$ .

**Example 4.4.1.** Estimating a Bernoulli parameter.

Let  $X_1, \dots, X_n$  be iid with

$$\begin{aligned} P_\theta(X_i = x) &= \theta^x(1 - \theta)^{1-x}, \quad x = 0, 1 \\ A &= [0, 1] = \Theta, \quad L(\theta, a) = (\theta - a)^2. \end{aligned}$$

This defines a statistical decision problem, which is the problem of estimating  $\theta$  under squared-error loss. Consider a class of estimators (ie, functions of the data taking values in  $\Theta$ ) denoted by  $d_\gamma(\mathbf{x}) = \gamma\bar{x}$ ,  $0 \leq \gamma \leq 1$ , where  $\bar{x} = n^{-1} \sum_{i=1}^n x_i$ . The usual estimator  $d_1(\mathbf{x}) = \bar{x}$  belongs to this class. Is there an estimator  $d_{\gamma^*}(\mathbf{x})$  in this class which is better than all others?

*Solution.* Note that

$$\begin{aligned} R(\theta, d_\gamma) &= E_\theta[(\gamma\bar{X} - \theta)^2] = \text{Var}_\theta[\gamma\bar{X}] + \{E_\theta[\gamma\bar{X}] - \theta\}^2 \\ &= \frac{\gamma^2\theta(1-\theta)}{n} + (\gamma-1)^2\theta^2. \end{aligned}$$

So  $R(\theta, d_\gamma)/R(\theta, d_1) = \gamma^2 + n\theta(1-\gamma)^2/(1-\theta)$ , which is  $>$ ,  $=$ , or  $< 1$  according as  $\theta >$ ,  $=$ , or  $< (1+\gamma)/[n(1-\gamma) + (1+\gamma)]$ . For example, taking  $\gamma = 3/4$ ,  $d_\gamma$  is better or worse than  $d_1$  according as  $\theta <$  or  $> 7/(n+7)$ . This shows that there is no clear winner in the class of  $d_\gamma$ ,  $0 \leq \gamma \leq 1$ .

**Example 4.4.2.** Choosing between two values of  $\theta$ .

Let  $X \sim P_\theta$ ,  $\Theta = \{\theta_0, \theta_1\}$ ,  $A = \{a_0, a_1\}$  where  $a_i$  is to choose  $\theta_i$  as the true value of  $\theta$ , and let  $L(\theta_0, a_0) = L(\theta_1, a_1) = 0$  and  $L(\theta_0, a_1) = L(\theta_1, a_0) = 1$ . A typical behavioral decision rule is described by a function  $\varphi: \mathfrak{X} \rightarrow [0, 1]$ , so that

$$\delta_\varphi(a_1|x) = \varphi(x) \text{ and } \delta_\varphi(a_0|x) = 1 - \varphi(x),$$

that is,  $\delta_\varphi$  chooses actions  $a_1, a_0$  with probabilities  $\varphi(x)$  and  $1 - \varphi(x)$ , respectively. Is there a best decision rule in this class?

*Solution.* The risk function of  $\delta_\varphi$  is given by

$$\begin{aligned} R(\theta_0, \delta_\varphi) &= \int_{\mathfrak{X}} [L(\theta_0, a_0)\{1 - \varphi(x)\} + L(\theta_0, a_1)\varphi(x)]f(x, \theta_0) dx \\ &= E_{\theta_0}[\varphi(X)], \\ R(\theta_1, \delta_\varphi) &= \int_{\mathfrak{X}} [L(\theta_1, a_0)\{1 - \varphi(x)\} + L(\theta_1, a_1)\varphi(x)]f(x, \theta_1) dx \\ &= 1 - E_{\theta_1}[\varphi(X)]. \end{aligned}$$

The *risk set*  $S = \{(R(\theta_0, \delta_\varphi), R(\theta_1, \delta_\varphi)), \varphi: \mathfrak{X} \rightarrow [0, 1]\}$  has the following features:

- (i)  $(0, 1) \in S$  corresponding to  $\varphi(x) \equiv 0$ .
- (ii)  $(1, 0) \in S$  corresponding to  $\varphi(x) \equiv 1$ .
- (iii)  $S$  is convex, because for any  $\varphi_1, \varphi_2, \lambda\varphi_1 + (1-\lambda)\varphi_2$  for  $0 \leq \lambda \leq 1$ , has the property:

$$\begin{aligned} &(R(\theta_0, \delta_{\lambda\varphi_1 + (1-\lambda)\varphi_2}), R(\theta_1, \delta_{\lambda\varphi_1 + (1-\lambda)\varphi_2})) \\ &= \lambda(R(\theta_0, \delta_{\varphi_1}), R(\theta_1, \delta_{\varphi_1})) + (1-\lambda)(R(\theta_0, \delta_{\varphi_2}), R(\theta_1, \delta_{\varphi_2})). \end{aligned}$$

- (iv)  $(\alpha, \beta) \in S$  implies  $(1-\alpha, 1-\beta) \in S$  (ie,  $S$  is symmetric about  $(1/2, 1/2)$ ). (To see this, consider  $\psi = 1 - \varphi$  for any  $\varphi$ .)
- (v)  $S$  can be shown to be closed. We omit the proof.

Consider any two points on the lower boundary of  $S$ . Among the  $\varphi$ 's corresponding to two such points, there is no clear winner. However, for any  $s \in S$ , the point  $s^*$  on the lower boundary which is on the vertical strip through  $s$  is decidedly better than  $s$ . For this reason, it is enough to restrict our attention to the  $\varphi$ 's corresponding to the points on the lower boundary of  $S$ .

## 4.5 Sufficient Statistics

Consider a problem of statistical inference in the framework of a parametric family with pdf or pmf  $\{f(\mathbf{x}, \theta), \theta \in \Theta\}$  based on a random sample  $(X_1, \dots, X_n)$ . The unknown  $\theta \in \Theta$ , which we want to estimate or test a hypothesis about, is of a finite dimension, often of only one or two dimensions, while the number of observations in the sample,  $n$ , may be quite large. Do we really need to carry all the  $n$  observations in the data in our search for a good estimate or test or whatever else to decide about?

**Definition 4.5.1.** A function  $T: \mathfrak{X}^n \rightarrow \mathbb{R}$  of the sample observations is said to be a statistic if

$$\{(x_1, \dots, x_n) : T(x_1, \dots, x_n) \leq a\} \in \mathbb{B}^n \quad \text{for all } a \in \mathbb{R},$$

where  $\mathbb{B}^n$  is the class of all Borel sets in  $\mathfrak{X}^n$ . A vector  $\mathbf{T} = (T_1, \dots, T_k)$  is said to be a  $k$ -dim statistic.

If we could extract all the relevant information about  $\theta$  in a *fixed-dimensional* statistic  $T = T(X_1, \dots, X_n)$ , then we can concentrate our efforts in search of a good procedure based on  $T$  without carrying the burden of the entire sample  $(X_1, \dots, X_n)$ . If we have such a statistic, then it can rightly be called a *Sufficient Statistic*.

**Definition 4.5.2.** A statistic  $T = T(X_1, \dots, X_n)$  is said to be sufficient for  $\theta \in \Theta$  in  $\mathbf{X}$  if the conditional distribution of  $\mathbf{X}$  given  $T$  is independent of  $\theta$ .

This means the following:

*Discrete Case.* Here the pmf of  $T$  is  $g(t, \theta) = \sum_{\{\mathbf{x}: T(\mathbf{x})=t\}} f(\mathbf{x}, \theta)$ , and the conditional distribution of  $\mathbf{X}$  given  $T = t$  is given by

$$h_\theta(\mathbf{x}|t) = \frac{P_\theta[\mathbf{X} = \mathbf{x}, T = t]}{P_\theta[T = t]} = \begin{cases} f(\mathbf{x}, \theta)/g(t, \theta) & \text{if } T(x_1, \dots, x_n) = t, \\ 0 & \text{otherwise} \end{cases}$$

for all  $t$  with  $g(t, \theta) > 0$ . Sufficiency requires  $h_\theta(\mathbf{x}|t) = h(\mathbf{x}|t)$ , independent of  $\theta$ .

*Continuous Case.* For simplicity of notations, suppose  $\theta \in \mathbb{R}$  and let  $T = T(X_1, \dots, X_n)$  be a 1-dim statistic. (Generalization to higher dimension is routine.) Suppose that  $U_i = U_i(X_1, \dots, X_n)$ ,  $1 \leq i \leq n - 1$  is any collection of  $n - 1$  other statistics such that  $\varphi: (X_1, \dots, X_n) \leftrightarrow (T, U_1, \dots, U_{n-1})$  is one-to-one with continuous first partials. Then the joint pdf of  $(T, U_1, \dots, U_{n-1})$  is

$$f(\varphi^{-1}(t, u_1, \dots, u_{n-1}), \theta) \left| J_\varphi(\varphi^{-1}(t, u_1, \dots, u_{n-1})) \right|^{-1}.$$

From this, the pdf of  $g(t, \theta)$  of  $T$  is obtained by integrating out  $u_1, \dots, u_{n-1}$  and the conditional pdf of  $(U_1, \dots, U_{n-1})$  given  $T = t$  is

$$h_\theta(u_1, \dots, u_{n-1}|t) = \frac{f(\varphi^{-1}(t, u_1, \dots, u_{n-1}), \theta) \left| J_\varphi(\varphi^{-1}(t, u_1, \dots, u_{n-1})) \right|^{-1}}{g(t, \theta)}.$$

Sufficiency of  $T$  requires  $h_\theta(u_1, \dots, u_{n-1}|t)$  to be independent of  $\theta$ . Since  $J_\varphi(\varphi^{-1}(t, u_1, \dots, u_{n-1}))$  does not involve  $\theta$  anyway, and since

$$f(\varphi^{-1}(T(\mathbf{x}), u_1(\mathbf{x}), \dots, u_{n-1}(\mathbf{x})), \theta) = f(\mathbf{x}, \theta),$$

$h_\theta(u_1, \dots, u_{n-1}|t)$  being independent of  $\theta$  is equivalent to  $f(\mathbf{x}, \theta)/g(T(\mathbf{x}), \theta)$  being independent of  $\theta$ . Thus in both the discrete and continuous case,  $T$  is sufficient for  $\theta$  in  $\mathbf{X}$  if  $f(\mathbf{x}, \theta)/g(T(\mathbf{x}), \theta)$  is independent of  $\theta$ .

If  $\mathbf{T}(\mathbf{x}) = (T_1(\mathbf{x}), \dots, T_k(\mathbf{x}))$  has the above property (using  $U_1, \dots, U_{n-k}$  to construct a one-to-one map), then  $(T_1, \dots, T_k)$  are jointly sufficient for  $\theta$ , usually when  $\theta = (\theta_1, \dots, \theta_k)$ .

Suppose  $T$  is sufficient for  $\theta$  in  $\mathbf{X}$  and let  $\psi$  be a real-valued function of  $\mathbf{X}$  with  $E_\theta[\psi(\mathbf{X})]$  finite. Let  $U_1, \dots, U_{n-1}$  be as described above, and  $\psi^*(T, U_1, \dots, U_{n-1}) = \psi \circ \varphi^{-1}(T, U_1, \dots, U_{n-1})$ . Then

$$\begin{aligned} E_\theta[\psi(\mathbf{X})|T = t] &= E_\theta[\psi^*(T, U_1, \dots, U_{n-1})|T = t] \\ &= E_\theta[\psi^*(t, U_1, \dots, U_{n-1})|T = t] \end{aligned}$$

is independent of  $\theta$ .

Now consider an arbitrary decision problem described by  $(\Theta, A, L)$  in the context of  $(\mathfrak{X}, \mathbb{B}, \{P_\theta, \theta \in \Theta\})$ . For simplicity of discussion, let  $\mathfrak{X}$  and  $A$  be countable and let  $\delta(\cdot|\mathbf{x})$  be an arbitrary behavioral rule. Suppose  $T$  is sufficient for  $\theta$  in  $\mathbf{X}$  and let  $\mathcal{T} = T(\mathfrak{X})$  and  $\mathfrak{X}_t = \{\mathbf{x} \in \mathfrak{X}: T(\mathbf{x}) = t\}$ . Then the risk of  $\delta$  is

$$\begin{aligned} R(\theta, \delta) &= \sum_{\mathbf{x} \in \mathfrak{X}} \left( \sum_{a \in A} L(\theta, a) \delta(a|\mathbf{x}) \right) f(\mathbf{x}, \theta) \\ &= \sum_{t \in \mathcal{T}} \sum_{\mathbf{x} \in \mathfrak{X}_t} \left( \sum_{a \in A} L(\theta, a) \delta(a|\mathbf{x}) \right) g(t, \theta) h(\mathbf{x}|t) \\ &= \sum_{t \in \mathcal{T}} \left[ \sum_{a \in A} L(\theta, a) \left( \sum_{\mathbf{x} \in \mathfrak{X}_t} \delta(a|\mathbf{x}) h(\mathbf{x}, t) \right) \right] g(t, \theta) \\ &= \sum_{t \in \mathcal{T}} \left( \sum_{a \in A} L(\theta, a) \delta^*(a|t) \right) g(t, \theta) \\ &= R(\theta, \delta^*), \end{aligned}$$

where  $\delta^*$  is defined by

$$\delta^*(a|\mathbf{x}) = \delta^*(a|t) = \sum_{\mathbf{x} \in \mathfrak{X}_t} \delta(a|\mathbf{x}) h(\mathbf{x}, t) \quad \text{for all } \mathbf{x} \in \mathfrak{X}_t,$$

which depends on  $\mathbf{x}$  only through  $T(\mathbf{x})$  and so we write  $\delta^*(a|\mathbf{x}) = \delta^*(a|t)$ . In other words,  $\delta^*(a|t) = E[\delta(a|X)|T = t]$ . We thus have for arbitrary  $\delta$ , an equivalent decision rule  $\delta^*$  which uses the data  $\mathbf{x}$  only through the summary provided by the statistic  $T$ . This justifies the term sufficient statistic. In any decision problem, a sufficient statistic is used for data reduction without losing anything that could be achieved by the full data. We now state and prove the key theorem about sufficient statistic.

**Theorem 4.5.1** (Factorization Theorem). *A statistic  $T = T(\mathbf{X})$  is sufficient for  $\theta$  in  $\mathbf{X}$  iff there exist functions  $g(t, \theta)$  and  $h(\mathbf{x})$  so that*

$$f(\mathbf{x}, \theta) = g(T(\mathbf{x}), \theta) h(\mathbf{x}) \quad \text{for all } \mathbf{x}, \theta,$$

where for every fixed  $t = T(\mathbf{x})$ , the function  $h(\mathbf{x})$  is independent of  $\theta$ .

The proof follows from our discussion above.

**Example 4.5.1.** In [Section 2.2.11](#) we defined an exponential family of distributions of  $\mathbf{X}$  without specifying  $\mathfrak{X}$  in which  $\mathbf{X}$  takes its values. We now let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a  $k$ -parameter regular exponential family. Then

$$f_{\mathbf{X}}(\mathbf{x}, \theta) = \{c(\theta)\}^n \exp \left[ \sum_{j=1}^k Q_j(\theta) \sum_{i=1}^n T_j(x_i) \right] \prod_{i=1}^n r(x_i).$$

Thus  $(T_1, \dots, T_k) = \sum_{i=1}^n (T_1(x_i), \dots, T_k(x_i))$  are jointly sufficient for  $\theta$  in  $\mathbf{X}$ .

## 4.6 Optimal Decision Rules

**Definition 4.6.1.** A decision rule  $\delta_1$  is said to be (i) as good as  $\delta_2$ , (ii) better than  $\delta_2$ , or (iii) equivalent to  $\delta_2$ , according as

- (i)  $R(\theta, \delta_1) \leq R(\theta, \delta_2)$  for all  $\theta \in \Theta$ ,
- (ii)  $R(\theta, \delta_1) < R(\theta, \delta_2)$  for all  $\theta \in \Theta$ , with strict inequality for some  $\theta$ ,
- (iii)  $R(\theta, \delta_1) = R(\theta, \delta_2)$  for all  $\theta \in \Theta$ .

If there exists a decision rule  $\delta_0$  which is as good as any other  $\delta$ , then of course  $\delta_0$  is optimal, but typically, such a  $\delta_0$  does not exist as seen in [Examples 4.4.1](#) and [4.4.2](#). The concept of optimality therefore needs adjustment. Two general approaches are taken for this purpose:

- (a) restricting the class of decision rule to choose from, or
- (b) ordering the decision rules in a less stringent manner.

### 4.6.1 Restrictions Used in the Estimation Problem

#### *Unbiasedness*

Restrict attention to only those estimators  $\hat{\theta} = d(X)$  of  $\theta$ , which satisfies

$$E_{\theta}[\hat{\theta}] = \theta \quad \text{for all } \theta \in \Theta.$$

In the example of estimating a Bernoulli parameter, all  $d_\gamma$  with  $\gamma \neq 1$  are now ruled out. In a large class of estimation problems with squared-error loss, there exists a best estimator in the class of unbiased estimators. These are the uniformly minimum variance unbiased estimators (UMVUE).

#### *Equivariance*

Suppose that  $f(x, \theta) = g(x - \theta)$  where  $g$  is a known pdf and we want to estimate  $\theta$ , which is called a location parameter, subject to  $L(\theta, a) = (a - \theta)^2$ . If  $X_1, \dots, X_n$  are iid with pdf  $f(\cdot, \theta)$ , then  $X_1 + c, \dots, X_n + c$  are iid with pdf  $f(\cdot, \theta + c)$ . Moreover,  $L(\theta + c, a + c) = L(\theta, a) = (a - \theta)^2$ . Thus the problem of estimating  $\theta + c$  from  $X_1 + c, \dots, X_n + c$  is the same as that of estimating  $\theta$  from  $X_1, \dots, X_n$ ; that is, the problem of estimating a location parameter under squared-error loss is *invariant under location*, that is, under the transformations

$$g_c: \mathfrak{X} \rightarrow \mathfrak{X} \text{ and its corresponding } \bar{g}_c: \Theta \rightarrow \Theta, \quad c \in \mathbb{R},$$

defined by  $g_c(x) = x + c$  and  $\bar{g}_c(\theta) = \theta + c$ . It is therefore reasonable to restrict our estimator by the requirement:  $d(x_1 + c, \dots, x_n + c) = d(x_1, \dots, x_n) + c$  for all  $x, c$ . Such estimators are called *equivariant under location*. Among equivariant estimators in this, and many more general invariant problems, there exists a best estimator, known as minimum risk equivariant estimator (MRE).

Optimal estimators under these restrictions will be discussed in [Chapter 5](#).

### 4.6.2 Restriction Used in the Two-Decision Problem

As mentioned earlier, the two-decision problem is treated as a problem of testing a hypothesis  $H_0: \theta \in \Theta_0$  against  $H_1: \theta \in \Theta_1$  by introducing an asymmetry between the two hypotheses. We are *inclined to accept*  $H_0$  unless there is *significant evidence provided by the data to reject  $H_0$  in favor of  $H_1$* . For this reason, we call  $H_0$  the *null hypothesis* and the error in rejecting  $H_0$  when it is true is called a *Type I error*, while  $H_1$  is called the *alternative hypothesis* and the error in accepting  $H_0$  when  $H_0$  is not true (ie,  $H_1$  is true) is called a *Type II error*.

We illustrate the asymmetric nature of  $H_0$  and  $H_1$  by the following example. Suppose for a certain disease there is a drug on the market with an effectiveness of  $p_0 = 0.40$  established over a long period of time. Now a new drug which claims to be better has been shown to be effective in  $\hat{p}_1 = 0.45$  in a clinical study based on  $n = 100$  cases. To test the validity of this claim, we have to choose between the actions:

$$a_0 = \{p_1 \leq 0.40\} \text{ and } a_1 = \{p_1 > 0.40\},$$

where  $p_1$  = the unknown effectiveness of the new drug. Here  $H_0: p_0 = 0.40$  is based on a long record, while  $H_1: p_1 > p_0 = 0.40$  needs to be substantiated by significant (ie, overwhelming) evidence provided by the data in favor of  $H_1$ . This is why  $H_0$  is called the *null hypothesis*, giving it a special role, while  $H_1$  is called an *alternative hypothesis* which is to be accepted only if the evidence provided by the data in its favor is significant enough. Therefore, we want to control the probability,  $P_{H_0}[\text{Reject } H_0] \leq \alpha$  (preassigned) and with that restriction, want to maximize  $P_{H_1}[\text{Reject } H_0]$ . We now introduce the restrictions in the two-decision problem.

### *Prescribed Type I Error Probability*

Recall the discussion about the lower boundary of the risk set in [Example 4.4.2](#). Here if we require

$$R(\theta_0, \delta_\varphi) = E_{\theta_0}[\varphi(X)] = P_{\theta_0}[\text{Decide } \theta = \theta_1] = \alpha,$$

where  $0 < \alpha < 1$  is given, then we are restricted to a vertical strip in the risk set, and in this restricted set, the  $\varphi^*$  which corresponds to the point  $s^*$  on the lower boundary in that vertical strip is the best (ie, this  $\varphi^*$  minimizes the Type II Error Probability  $R(\theta_1, \delta_\varphi) = 1 - E_{\theta_1}[\varphi(X)] = P_{\theta_1}[\text{Decide } \theta = \theta_0]$  subject to the requirement that the Type I Error Probability  $R(\theta_0, \delta_\varphi) = P_{\theta_0}[\text{Decide } \theta = \theta_1] = \alpha$ ).

### *Unbiased Tests*

The problem of testing  $H_0: \theta \in \Theta_0$  against  $H_1: \theta \in \Theta_1$  where  $\Theta_0$  and  $\Theta_1$  are disjoint sets in  $\Theta$  is more complicated than testing  $H_0: \theta = \theta_0$  against  $H_1: \theta = \theta_1$ . Here we restrict to those  $\delta_\varphi$  for which

- (i)  $\sup_{\theta \in \Theta_0} E_\theta[\varphi(X)] = \alpha$  (these are tests of level  $\alpha$ ),
- (ii)  $E_\theta[\varphi(X)] \geq \alpha$  for all  $\theta \in \Theta_1$  (these are unbiased tests of level  $\alpha$ ),

and then among tests satisfying (i) and (ii) search for  $\varphi^*$  in this class for which

- (iii)  $E_\theta[\varphi^*(X)] \geq E_\theta[\varphi(X)]$ , for all  $\varphi$ , and for all  $\theta \in \Theta_1$ .

In many situations such a  $\varphi^*$  exists and is called the *Uniformly Most Powerful (UMP) Unbiased* level  $\alpha$  test.

**Note.** If a test  $\varphi$  does not satisfy condition (ii) for unbiasedness, then  $E_{\theta_1}[\varphi(X)] = \alpha - \varepsilon$  for some  $\theta_1 \in \Theta_1$  and some  $\varepsilon > 0$ , while  $E_{\theta_0}[\varphi(X)] > \alpha - \varepsilon$  for some  $\theta_0 \in \Theta_0$ , so that  $\varphi$  rejects  $H_0$  when  $\theta = \theta_0$  and  $H_0$  is true with a larger probability than when  $\theta = \theta_1$  and  $H_0$  is not true. The condition of unbiasedness does not allow such undesirable decisions.

Optimum tests under these restrictions will be discussed in [Chapter 6](#).

### 4.6.3 Suitable Ordering of Decision Rules

The stringent ordering by the entire risk function can be replaced by invoking other principles of ordering. The following are two such important principles.

### The Bayes Principle

Let  $\tau$  denote a probability distribution on the parameter space  $\Theta$ . Actually, we should introduce a suitable class of events  $\mathcal{C}$  in  $\Theta$  at this point, on which  $\tau$  is defined, but when  $\Theta = \mathbb{R}^k$  or a subset thereof (as in many situations), we shall simply use the events  $\mathbb{B}$  (introduced earlier) for this purpose. Moreover, with a slight abuse of notation, we shall let  $\tau$  denote the pdf of this distribution and define the Bayes risk of a behavioral decision rule  $\delta$  with respect to  $\tau$  as:

$$r(\tau, \delta) = \int_{\Theta} R(\theta, \delta) \tau(\theta) d\theta$$

with  $R(\theta, \delta)$  as defined in [Section 4.4](#).

**Definition 4.6.2.** The probability distribution  $\tau$  is called the *prior distribution of  $\theta$* . A decision rule  $\delta_0$  is said to be *Bayes* with respect to the prior distribution  $\tau$  if

$$r(\tau, \delta_0) = \inf_{\delta} r(\tau, \delta).$$

There may be many Bayes rules with respect to a  $\tau$ .

Sometimes it is useful to work with a function  $\tau$  on  $\Theta$  which is *not a pdf* and define  $r(\tau, \delta)$  for such a  $\tau$  in a formal way, and then minimize it with respect to  $\delta$ .

**Definition 4.6.3.** A decision rule  $\delta_0$  is said to be

(a) generalized Bayes if there exists  $\tau(\theta) \geq 0$  but  $\int_{\Theta} \tau(\theta) d\theta = \infty$  (ie,  $\tau$  is not a pdf), and

$$\int_{\Theta} R(\theta, \delta_0) \tau(\theta) d\theta \leq \int_{\Theta} R(\theta, \delta) \tau(\theta) d\theta \quad \text{for all } \delta,$$

(b)  $\varepsilon$ -Bayes with respect to  $\tau$  if  $r(\tau, \delta_0) \leq \inf_{\delta} r(\tau, \delta) + \varepsilon$ ,  $\varepsilon > 0$ , and

(c) extended Bayes if  $\delta_0$  is  $\varepsilon$ -Bayes for every  $\varepsilon > 0$  (ie, for every  $\varepsilon > 0$ , there exists a prior distribution  $\tau_{\varepsilon}$  such that  $\delta_0$  is  $\varepsilon$ -Bayes with respect to  $\tau_{\varepsilon}$ ).

### The Minimax Principle

**Definition 4.6.4.** A decision rule  $\delta_0$  is said to be *minimax* if

$$\sup_{\theta} R(\theta, \delta_0) = \inf_{\delta} \sup_{\theta} R(\theta, \delta).$$

There may be many minimax rules.

The minimax principle summarizes the performance of each  $\delta$  by  $\sup_{\theta} R(\theta, \delta)$  (ie, judges a rule  $\delta$  by its worst performance). A minimax rule may, therefore, have quite weak overall performance.

### Admissibility

Finally, we introduce another property of a decision rule which requires that it cannot be improved upon.

**Definition 4.6.5.** A decision rule  $\delta_0$  is said to be *inadmissible* if there is a rule  $\delta_1$  which is better than  $\delta_0$  (ie,  $R(\theta, \delta_1) \leq R(\theta, \delta_0)$  for all  $\theta \in \Theta$  with strict inequality for some  $\theta \in \Theta$ ).

A decision rule  $\delta_0$  is said to be *admissible* if there is no decision rule that is better than  $\delta_0$ , that is, for every rule  $\delta$ ,  $R(\theta_1, \delta) < R(\theta_1, \delta_0)$  for some  $\theta_1 \in \Theta$  implies that there exists  $\theta_2 \in \Theta$  such that  $R(\theta_2, \delta) > R(\theta_2, \delta_0)$ .

Admissibility is in no sense an indicator of high performance of a decision rule. It merely says that there is none that is uniformly better. On the other hand, inadmissibility indicates that there is a uniformly better rule that we may want to search for. However, there are rules with many nice properties which are inadmissible. A prime example is the sample mean vector  $\bar{\mathbf{X}}_n$  as an estimator of the mean vector of  $N_p(\boldsymbol{\mu}, I)$  with the loss function

$$L(\boldsymbol{\mu}, \mathbf{a}) = \|\mathbf{a} - \boldsymbol{\mu}\|^2 = (\mathbf{a} - \boldsymbol{\mu})^T (\mathbf{a} - \boldsymbol{\mu})$$

for  $p \geq 3$ , which is strictly improved upon by the shrinkage estimator

$$\delta_c(\mathbf{X}) = \bar{\mathbf{X}}_n - \frac{p-2}{\|\bar{\mathbf{X}}_n - \mathbf{c}\|^2} (\bar{\mathbf{X}}_n - \mathbf{c}), \quad \mathbf{c} \in \mathbb{R}^p$$

due to Stein [11]. Unfortunately,  $\delta_c$  is also inadmissible, because it can be improved upon by (see [3, p. 302–3])

$$\delta_c^+(\mathbf{X}) = \bar{\mathbf{X}}_n - \min\left\{1, \frac{p-2}{\|\bar{\mathbf{X}}_n - \mathbf{c}\|^2}\right\} (\bar{\mathbf{X}}_n - \mathbf{c}).$$

But  $\delta_c^+$  is also inadmissible.

#### 4.6.4 Finding Bayes Rules: Prior to Posterior

The idea in the Bayes Principle is to judge a  $\delta$  by its average performance with weights assigned by  $\tau$ , which reflects the likelihood of various values of  $\theta$  in the statistician's assessment based on prior experience (before obtaining the data). The key step in the calculation of the Bayes rules is to incorporate the data  $X = x$  in the prior distribution  $\tau$  to obtain the *posterior distribution of  $\theta$* , which is the continuous analog of the Bayes Formula in [Chapter 1, Proposition 1.6.1](#). Write  $f(x, \theta) = f(x|\theta)$ , treating  $f(x, \theta)$  as the conditional pdf of  $X$  given  $\theta$ . Then proceeding as in [Section 1.10](#), we have

$$\begin{aligned} r(\tau, d) &= \int_{\Theta} \left[ \int_{\mathfrak{X}} L(\theta, d(x)) f(x|\theta) dx \right] \tau(\theta) d\theta \\ &= \int_{\mathfrak{X}} \left[ \int_{\Theta} L(\theta, d(x)) \frac{f(x|\theta)\tau(\theta)}{\int_{\Theta} f(x|\theta)\tau(\theta) d\theta} d\theta \right] \left( \int_{\Theta} f(x|\theta)\tau(\theta) d\theta \right) dx \\ &= \int_{\mathfrak{X}} \left[ \int_{\Theta} L(\theta, d(x)) g(\theta|x) d\theta \right] f(x) dx \\ &= \int_{\mathfrak{X}} E[L(\theta, d(x))|X = x] f(x) dx, \end{aligned}$$

where  $f(x) = \int_{\Theta} f(x|\theta)\tau(\theta) d\theta$  is the *marginal* pdf of  $X$ , and  $g(\theta|x) = [f(x|\theta)\tau(\theta)] / \int_{\Theta} f(x|\theta)\tau(\theta) d\theta$  is the *posterior* pdf of  $\theta$  given  $X = x$ . The minimization of  $r(\tau, d)$  with respect to  $d$  is achieved by choosing  $d(x) = d^*(x) \in A$  for each  $x$  so that

$$E[L(\theta, d^*(X))|X = x] \leq E[L(\theta, a)|X = x] \quad \text{for all } a \in A.$$

Interchanging the order of integration over  $\Theta$  and  $\mathfrak{X}$  is justified if  $L(\theta, a)$  is bounded below. The case when  $X$  is discrete is treated in exactly the same way, replacing  $\int_{\mathfrak{X}} dx$  by  $\sum_{\mathfrak{X}}$ .

**Example 4.6.1.** Estimating the mean of  $N(\theta, \sigma^2)$  under squared-error loss.

Our data consists of  $X_1, \dots, X_n$  iid as  $N(\theta, \sigma^2)$  with  $\sigma^2$  known, and  $\Theta = A = \mathbb{R}$ ,  $L(\theta, a) = (a - \theta)^2$ . Since  $\bar{X}_n$  is sufficient for  $\theta$  in  $\mathbf{X}$ , it is enough to restrict attention to decision rules based on  $\bar{X}_n$  having pdf

$$f(\bar{x}_n, \theta) = \frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \exp\left[-\frac{n}{2\sigma^2}(\bar{x}_n - \theta)^2\right], \quad \text{ie, } \bar{X}_n \sim N\left(\theta, \frac{\sigma^2}{n}\right).$$

First consider the prior  $\tau$ :

$$\tau(\theta) = \frac{1}{\sqrt{2\pi}\gamma} \exp\left[-\frac{1}{2\gamma^2}(\theta - \mu)^2\right],$$

that is,  $\theta \sim N(\mu, \gamma^2)$ , and  $\bar{X}_n|\theta \sim N\left(\theta, \frac{\sigma^2}{n}\right)$ . Using the properties of bivariate normal distribution (see [Section 2.2.10](#)), we have

$$\theta|\bar{X}_n \sim N\left(\frac{w_1\mu + w_2\bar{X}_n}{w_1 + w_2}, \frac{1}{w_1 + w_2}\right)$$

as the posterior distribution, where  $w_1 = 1/\gamma^2$  and  $w_2 = n/\sigma^2$ , that is,  $E[\theta|\bar{X}_n]$  is the weighted average of the mean  $\mu$  of the prior and the sample mean  $\bar{X}_n$  with weights inversely proportional to the variance of the prior distribution for  $\theta$  and the conditional variance of  $\bar{X}_n$  given  $\theta$  for  $\bar{X}_n$ .

*Remark 4.6.1.*

1. This effect of the prior distribution on  $\bar{X}_n$  is called: “a shrinkage of the sample mean toward the prior mean.”
2. On the other hand, we can also interpret this as “overcoming a prior belief by observed data.”
3. Note that the weight  $1/\gamma^2$  attached to the prior remains fixed, but the weight  $n/\sigma^2$  attached to  $\bar{X}_n$  keep increasing with  $n$  and  $1/\sigma^2$  determines the rate of increase. Thus a bad prior cannot hurt much once the sample size gets large.

Finally, the Bayes estimator of  $\theta$  under squared-error loss, using this prior is:

$$E[\theta|\bar{X}_n] = \frac{n\gamma^2\bar{X}_n + \sigma^2\mu}{n\gamma^2 + \sigma^2}.$$

For the special case of  $\sigma^2 = 1$ ,  $n = 1$ , and  $\mu = 0$ ,  $E[\theta|X] = \frac{\gamma^2}{1+\gamma^2}X$ .

### 4.6.5 Solving for Minimax Rules

In this section we discuss two methods of finding minimax rules.

**Definition 4.6.6.** A prior distribution  $\tau_0$  is said to be *least favorable* if

$$\inf_{\delta} r(\tau_0, \delta) = \sup_{\tau} \inf_{\delta} r(\tau, \delta).$$

Next note that

$$\underline{V} = \sup_{\tau} \inf_{\delta} r(\tau, \delta) \leq \inf_{\delta} \sup_{\tau} r(\tau, \delta) = \overline{V},$$

where  $\underline{V}$  and  $\overline{V}$  are, respectively, the lower and upper value of the statistical decision problem which is viewed as a game between nature who chooses  $\theta$  randomly according to the distribution  $\tau$  and the statistician who chooses a decision rule  $\delta$ , resulting in a payoff by the statistician of the quantity  $r(\tau, \delta)$ . Also note that

$$\sup_{\theta} R(\theta, \delta) = \sup_{\tau} r(\tau, \delta) \quad \text{for all } \delta.$$

We say that the game of a statistical decision problem has a value if  $\underline{V} = \overline{V}$ .

We now describe the *first method* for finding a minimax rule which involves guessing a  $\tau_0$  as least favorable, finding a rule  $\delta_0$  which is Bayes with respect to  $\tau_0$  and checking that  $\delta_0$  is indeed minimax. The actual procedure is described in [Theorem 4.6.1](#).

In many situations, our guess of the least favorable  $\tau_0$  is not a pdf (ie,  $\int \tau(\theta) d\theta = \infty$ ), in which case  $\delta_0$ , which is formally Bayes with respect to  $\tau_0$ , is not really a Bayes rule. This needs a modification of the above method, which is described in [Theorem 4.6.2](#).

**Theorem 4.6.1.** *If  $\delta_0$  is Bayes with respect to  $\tau_0$  and*

$$R(\theta, \delta_0) \leq r(\tau_0, \delta_0) \quad \text{for all } \theta,$$

*then  $\delta_0$  is minimax and  $\tau_0$  is least favorable.*

*Proof.* Note that

$$\begin{aligned} \overline{V} &= \inf_{\delta} \sup_{\tau} r(\tau, \delta) = \inf_{\delta} \sup_{\theta} R(\theta, \delta) \leq \sup_{\theta} R(\theta, \delta_0) \leq r(\tau_0, \delta_0) \\ &= \inf_{\delta} r(\tau_0, \delta) \leq \sup_{\tau} \inf_{\delta} r(\tau, \delta) = \underline{V}. \end{aligned}$$

Therefore, all inequalities are equalities and therefore,

$$\inf_{\delta} r(\tau_0, \delta) = \sup_{\tau} \inf_{\delta} r(\tau, \delta),$$

that is,  $\tau_0$  is least favorable, and

$$\sup_{\theta} R(\theta, \delta_0) = \inf_{\delta} \sup_{\theta} R(\theta, \delta),$$

which proves that  $\delta_0$  is minimax.  $\square$

**Theorem 4.6.2.** *Let  $\{\tau_n\}$  be a sequence of distributions of  $\theta$  and let  $\delta_n$  be Bayes with respect to  $\tau_n$ . If*

$$R(\theta, \delta_0) \leq \overline{\lim}_{n \rightarrow \infty} r(\tau_n, \delta_n) \quad \text{for all } \theta,$$

*then  $\delta_0$  is minimax.*

*Proof.* For arbitrary  $\delta$  and for all  $n$ ,

$$\int R(\theta, \delta) \tau_n(\theta) d\theta = r(\tau_n, \delta) \geq r(\tau_n, \delta_n),$$

implying  $\sup_{\theta} R(\theta, \delta) \geq r(\tau_n, \delta_n)$ . Therefore, for all  $\delta$ ,

$$\sup_{\theta} R(\theta, \delta) \geq \overline{\lim}_n r(\tau_n, \delta_n) \geq \sup_{\theta} R(\theta, \delta_0),$$

proving that  $\delta_0$  is minimax.  $\square$

*Remark 4.6.2.* Often  $c = \lim_n r(\tau_n, \delta_n)$  exists, in which case we simply have to check that  $R(\theta, \delta_0) \leq c$  for all  $\theta$ .

**Definition 4.6.7.** A decision rule is an equalizer rule if  $R(\theta, \delta_0)$  is constant for all  $\theta$ .

The second method of finding a minimax rule is to look for an equalizer rule that is Bayes or extended Bayes.

**Theorem 4.6.3.** *If an equalizer rule is extended Bayes, then it is minimax.*

*Proof.* Since  $R(\theta, \delta_0) = c$  (constant) for all  $\theta$ ,  $r(\tau, \delta_0) = c$  for all  $\tau$ , and since  $\delta_0$  is extended Bayes, there exists a sequence  $\{\tau_n\}$  with Bayes rules  $\{\delta_n\}$  such that

$$c = r(\tau_n, \delta_0) \leq r(\tau_n, \delta_n) + 1/n \quad \text{for all } n.$$

Hence for all  $\theta$ ,

$$R(\theta, \delta_0) = c \leq \overline{\lim}_n r(\tau_n, \delta_n).$$

It now follows from [Theorem 4.6.2](#) that  $\delta_0$  is minimax.  $\square$

#### 4.6.6 Conditions for Admissibility

The simplest (and obvious) condition for admissibility is the following.

**Theorem 4.6.4.** *If  $\delta_0$  is unique Bayes with respect to some  $\tau$ , then  $\delta_0$  is admissible.*

*Proof.* Since  $\delta_0$  is unique Bayes with respect to  $\tau$ ,  $r(\tau, \delta_0) < r(\tau, \delta_1)$  for all  $\delta_1 \neq \delta_0$ . Hence there exists  $\theta$  for any  $\delta_1 \neq \delta_0$  such that  $R(\theta, \delta_0) < R(\theta, \delta_1)$ .  $\square$

The requirement of uniqueness in the above theorem can be dispensed with if the risk function  $R(\theta, \delta)$  for every  $\delta$  is continuous in  $\theta$  and if the  $\tau$  with respect to which  $\delta_0$  is Bayes has plenty of support.

**Theorem 4.6.5.** *Suppose  $\Theta$  is an interval in  $\mathbb{R}$ , the risk function  $R(\theta, \delta)$  for every decision rule  $\delta$  is continuous in  $\theta$  and  $\tau$  is a probability distribution on  $\Theta$  such that  $\tau(I) > 0$  for all nondegenerate intervals  $I \subset \Theta$ . Then a Bayes rule  $\delta_0$  with respect to  $\tau$  is admissible.*

*Proof.* If  $\delta_0$  is inadmissible, then there is a rule  $\delta_1$  such that  $R(\theta, \delta_1) \leq R(\theta, \delta_0)$  for all  $\theta \in \Theta$ , and  $R(\theta_0, \delta_1) < R(\theta_0, \delta_0)$  for some  $\theta_0 \in \Theta$ . Since  $R(\theta, \delta_0)$  and  $R(\theta, \delta_1)$  are both continuous, there exist  $\varepsilon > 0$  and a neighborhood  $I$  of  $\theta_0$  such that  $R(\theta, \delta_0) - R(\theta, \delta_1) > \varepsilon$  for all  $\theta \in I$ . Hence

$$r(\tau, \delta_0) - r(\tau, \delta_1) \geq \int_I [R(\theta, \delta_0) - R(\theta, \delta_1)] \tau(\theta) d\theta > \varepsilon \tau(I) > 0,$$

contradicting the hypothesis of  $\delta_0$  being Bayes with respect to  $\tau$ .  $\square$

This theorem is now extended to cover some familiar situations where  $\delta_0$  is extended Bayes and has constant risk.

**Theorem 4.6.6.** *Suppose that a decision rule  $\delta_0$  satisfies the following conditions in relation to a sequence of probability distributions  $\{\tau_n\}$  on  $\Theta \subset \mathbb{R}$  with respect to which  $\{\delta_n\}$  are Bayes rules:*

- (i)  $R(\theta, \delta_0) = c$  (constant) for all  $\theta$ ,
- (ii)  $\lim_n r(\tau_n, \delta_n) = c$ ,
- (iii)  $R(\theta, \delta)$  is continuous in  $\theta$  for all  $\delta$ ,
- (iv)  $\tau_n(I) > 0$  for all nondegenerate intervals  $I \subset \Theta$ , and
- (v)  $\lim_n \tau_n(I)/[c - r(\tau_n, \delta_n)] = \infty$ .

Then  $\delta_0$  is admissible.

*Proof.* Note that

$$c - r(\tau_n, \delta_n) = r(\tau_n, \delta_0) - r(\tau_n, \delta_n) \geq 0 \quad \text{for all } n.$$

If equality holds for some  $n$ , then  $\delta_0$  is Bayes with respect to that  $\tau_n$  which makes it admissible by [Theorem 4.6.5](#). Therefore, assume that  $c - r(\tau_n, \delta_n) > 0$  for all  $n$  and suppose that  $\delta_0$  is inadmissible. Then there is a rule  $\delta'$  such that  $c = R(\theta, \delta_0) \geq R(\theta, \delta')$  for all  $\theta \in \Theta$ , and  $R(\theta_0, \delta_0) > R(\theta_0, \delta')$  for some  $\theta_0 \in \Theta$ . We can now find  $\varepsilon > 0$  and a neighborhood  $I$  of  $\theta_0$  so that  $c - r(\tau_n, \delta') > \varepsilon \tau_n(I)$  as in the proof of [Theorem 4.6.5](#). Thus

$$\frac{c - r(\tau_n, \delta')}{c - r(\tau_n, \delta_n)} > \frac{\varepsilon \tau_n(I)}{c - r(\tau_n, \delta_n)} > 1 \quad \text{for large } n,$$

since by Condition (v), the middle quantity in the last display tends to  $\infty$  as  $n \rightarrow \infty$ . This contradicts the hypothesis that  $\delta_n$  is Bayes with respect to  $\tau_n$ .  $\square$

[Theorems 4.6.5](#) and [4.6.6](#) depend on the key condition that the risk  $R(\theta, \delta)$  is continuous in  $\theta$  for all decision rules  $\delta$ , so the applicability of these theorems depends on the verification of this condition. The following theorem from Ferguson [1] serves that purpose for one-parameter exponential families.

**Theorem 4.6.7.** *Let  $\Theta$  be the real line. Suppose that*

- (a) *there exist nonnegative functions  $B_1(\theta_1, \theta_2)$  and  $B_2(\theta_1, \theta_2)$  bounded on compact sets of  $\Theta \times \Theta$  such that*

$$|L(\theta_2, a)| \leq B_1(\theta_1, \theta_2)|L(\theta_1, a)| + B_2(\theta_1, \theta_2), \quad \text{for all } a \in A;$$

- (b)  *$L(\theta, a)$  is continuous in  $\theta$  for each  $a \in A$ ; and*
- (c)  *$f(x, \theta) = c(\theta)h(x) \exp[Q(\theta)T(x)]$  where  $Q(\theta)$  is a continuous increasing function.*

Then for all nonrandomized decision rules  $d$ , the risk function  $R(\theta, d)$  is continuous in  $\theta$ .

*Proof.* See Ferguson [1], p. 139–40].  $\square$

Other methods of proving admissibility have been developed by Hodges and Lehmann [12], Karlin [13], and Stein [14].

**Example 4.6.2.** An admissible minimax estimator of a Bernoulli proportion under squared-error loss.

Let  $X$  be the number of successes in  $n$  independent Bernoulli trials with probability of success  $\theta$ . We want to estimate  $\theta$  under squared-error loss from the data  $X$ . Here  $\Theta = [0, 1] = A$  and  $L(\theta, a) = (a - \theta)^2$ .

*Solution.* We shall apply [Theorem 4.6.3](#), looking for an equalizer rule which is Bayes. This will give us a minimax rule. We look among Bayes rules with respect to  $\tau_{a,b} = \text{Beta}(a, b)$ , that is, Beta distributions with parameters  $a, b$  with pdf

$$\tau_{a,b}(\theta) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\theta^{a-1}(1-\theta)^{b-1}I_{(0,1)}(\theta).$$

The Bayes rule with respect to  $\tau_{a,b}$  is  $d_{a,b}(x) = (x + a)/(n + a + b)$  (the verification of which is left as an exercise), with risk function

$$\begin{aligned} R(\theta, d_{a,b}) &= E_\theta[d_{a,b}(X) - \theta]^2 = \text{Var}_\theta[d_{a,b}(X)] + \{E_\theta[d_{a,b}(X)] - \theta\}^2 \\ &= \frac{n\theta(1-\theta)}{(n+a+b)^2} + \left[ \frac{n\theta + a}{n+a+b} - \theta \right]^2 \\ &= \frac{[(a+b)^2 - n]\theta^2 - [2a(a+b) - n]\theta + a^2}{(n+a+b)^2}. \end{aligned}$$

To make this a constant in  $\theta$ , take  $(a+b)^2 = 2a(a+b) = n$ ; leading to  $a = b = \sqrt{n}/2$ . The resulting Bayes rule  $d_{\sqrt{n}/2, \sqrt{n}/2}(X) = (X + \sqrt{n}/2)/(n + \sqrt{n})$  is minimax.

Since this rule is unique Bayes with respect to  $\text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$ , it is admissible by [Theorem 4.6.4](#). Thus  $\hat{\theta} = (X + \sqrt{n}/2)/(n + \sqrt{n})$  is an admissible minimax estimator of  $\theta$ .

**Example 4.6.3.** Admissibility of  $\bar{X}_n$  as an estimator of mean  $\theta$  of  $N(\theta, \sigma^2)$  under squared-error loss.

Under squared-error loss, the Bayes estimator of  $\theta$  based on a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $N(\theta, \sigma^2)$ , with respect to  $\tau_k = N(0, k)$  is

$$d_k(\mathbf{X}) = \frac{nk}{nk + \sigma^2}\bar{X}_n, \quad \text{where } \bar{X}_n = n^{-1} \sum_{i=1}^n X_i.$$

This result is derived in [Section 4.6.4](#). The Bayes risk of  $d_k$  is  $r(\tau_k, d_k) = (k\sigma^2)/(nk + \sigma^2)$ , the proof of which is left an exercise. We shall use [Theorems 4.6.6](#) and [4.6.7](#) to show that the decision rule  $d_0(\mathbf{X}) = \bar{X}_n$  is admissible.

*Solution.* Since

$$R(\theta, d_0) = \sigma^2/n = \lim_{k \rightarrow \infty} r(\tau_k, d_k)$$

is constant, and  $\tau_k(I) > 0$  for all nondegenerate intervals, Conditions (i), (ii), and (iv) of [Theorem 4.6.6](#) are satisfied, while [Theorem 4.6.7](#) provides justification for Condition (iii) that  $R(\theta, d)$  is continuous for all  $d$ . So we only need to check Condition (v). Now

$$\frac{\tau_k([a, b])}{c - r(\tau_k, d_k)} = \frac{\Phi(b/\sqrt{k}) - \Phi(a/\sqrt{k})}{\sigma^2/n - (k\sigma^2)/(nk + \sigma^2)} = \frac{n(nk + \sigma^2)}{\sigma^4} [\Phi(b/\sqrt{k}) - \Phi(a/\sqrt{k})].$$

Since  $\Phi(b/\sqrt{k}) - \Phi(a/\sqrt{k}) = \frac{b-a}{\sqrt{2\pi k}}[1 + o(1)]$  as  $k \rightarrow \infty$  (see [Section 2.2.4](#)), the above tends to  $\infty$  as  $k \rightarrow \infty$ , justifying Condition (v). This proves the admissibility of  $\bar{X}_n$ .

## Exercises

- 4.1.** Suppose that  $X$  is distributed with a pdf  $f(x, \theta)$  where  $\theta$  is an unknown real parameter. Consider the two-decision problem of choosing between the hypotheses  $H_0: \theta \leq \theta_0$  and  $H_1: \theta > \theta_0$  (for a given  $\theta_0$ ) with the loss function

$$L(\theta, a_0) = (\theta - \theta_0)_+, \quad L(\theta, a_1) = (\theta_0 - \theta)_+,$$

where for any real number  $x$ ,  $x_+$  denotes  $\max(x, 0)$ , and  $a_i$  is the action to accept  $H_i$ ,  $i = 0, 1$ . Show that the Bayes rule with respect to a prior cdf  $G$  of  $\theta$  rejects  $H_0$  if and only if  $E_G[\theta|X = x] > \theta_0$ .

- 4.2.** Under 0 – 1 loss function in the problem of testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1$ , the risk set is

$$\begin{aligned} S &= \{(R(\theta_0, \phi), R(\theta_1, \phi)): \phi \text{ maps } \mathcal{X} \rightarrow [0, 1]\} \\ &= \{(E_{\theta_0}[\phi(X)], 1 - E_{\theta_1}[\phi(X)]): \phi \text{ maps } \mathcal{X} \rightarrow [0, 1]\}, \end{aligned}$$

where  $X$  denotes the data. Let  $f(x, \theta) = \binom{2}{x} \theta^x (1 - \theta)^{2-x}$ ,  $x = 0, 1, 2$ , and consider  $H_0: \theta = 1/2$  vs  $H_1: \theta = 2/3$ .

- (a) Plot the risk set  $S$ .
- (b) Find a minimax test.

- 4.3.** Let  $X$  follow a binomial distribution  $\text{Bin}(n, \theta)$ ,  $0 < \theta < 1$ .

- (a) Show that  $d_0(x) = x/n$  is a minimax estimator of  $\theta$  with constant risk  $1/n$  under the loss function  $L(\theta, a) = (\theta - a)^2 / \{\theta(1 - \theta)\}$ .
- (b) Show that  $d_0(x) = x/n$  is not minimax under the loss function  $L(\theta, a) = (\theta - a)^2$ . [Consider  $\delta_\varepsilon^*$  defined as:  $\delta_\varepsilon^*(x) = d_0(x) = x/n$  with probability  $1 - \varepsilon$  and  $\delta_\varepsilon^*(x) = d_1(x) \equiv 1/2$  with probability  $\varepsilon$ . Examine the risk function of  $\delta_\varepsilon^*$  with  $\varepsilon = (n+1)^{-1}$ .]
- (c) Show that  $d_0(x)$  is not a Bayes rule, but it is an extended Bayes rule (ie,  $\varepsilon$ -Bayes for every  $\varepsilon > 0$ ). [Try Beta priors.]
- (d) Show that  $d_0$  is Bayes with respect to  $\text{Unif}(0, 1)$  prior under  $L(\theta, a)$  given in (a).

- 4.4.** Let  $X$  follow the binomial  $\text{Bin}(n, \theta)$  distribution with pdf

$$f(x, \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}, \quad x = 0, 1, \dots, n,$$

where  $\theta \in \Theta = (0, 1)$  is to be estimated in the action space  $A = [0, 1]$  under the loss function  $L(\theta, a) = (a - \theta)^2$ .

- (a) Show that the estimator  $d_0(x) = x/n$  is admissible.  
 (b) If the parameter space is  $\Theta = [1/3, 2/3]$  and everything else is the same as above, then show that  $d_0(x) = x/n$  is inadmissible. [Hint: Think of a sensible estimator.]

**4.5.** Let  $X$  be an rv with pdf

$$f(x, \theta) = \theta^{-1} \exp(-x/\theta), \quad x > 0, \theta = 1 \text{ or } 2.$$

Consider the two-decision problem with action space  $A = \{a_1, a_2\}$  with  $a_i = i$  and the loss function  $L(\theta, a) = I(\theta \neq a)$ , based on a single observation  $X$ .

- (a) Let  $S$  denote the risk set consisting of points  $(R(1, \delta), R(2, \delta))$  corresponding to all behavioral rules  $\delta$ . Determine the function  $\beta(\alpha)$  describing the points  $(\alpha, \beta(\alpha))$  on the lower boundary of  $S$ .  
 (b) Find a minimax rule.  
 (c) Consider a prior  $g_\tau$  with probabilities  $\tau$  for  $\theta = 1$  and  $1 - \tau$  for  $\theta = 2$ . Show that the Bayes rule with respect to the prior  $g_\tau$  will always take action  $a_2$  (irrespective of the observed values of  $X$ ) if  $\tau < 1/3$ .  
**4.6.** Show that the sample mean  $\bar{X}$  is an admissible estimator of the mean  $\theta$  of a normal distribution  $N(\theta, \sigma^2)$  under the absolute error loss.  
**4.7.** Suppose we have one observation  $X$  from  $N(\theta, 1)$  on the basis of which we have to take action  $a_0$  to decide  $\theta = 0$  or  $a_1$  to decide  $\theta \neq 0$  subject the loss function

$$\begin{aligned} L(0, a_0) &= 0, & L(\theta, a_1) &= 0 \text{ for } \theta \neq 0, \\ L(0, a_1) &= 1, & L(\theta, a_0) &= 1 \text{ for } \theta \neq 0. \end{aligned}$$

Consider a prior distribution  $\tau$  which assigns probability  $p$  to  $\{\theta = 0\}$  and distributes the remaining probability  $1 - p$  on  $(-\infty, \infty)$  according to  $N(\mu, \sigma^2)$ , that is

$$\tau(\theta = 0) = p \text{ and } \tau[\theta \in B] = (1 - p) \int_B \sigma^{-1} \phi((\theta - \mu)/\sigma) d\theta \quad \text{if } 0 \notin B,$$

where  $\phi$  is the pdf of the standard normal distribution. Find the Bayes rule with respect to  $\tau$  based on a sample  $X$  of size 1. Give a common sense interpretation of the Bayes rule.

- 4.8.** Let  $X_1, \dots, X_n$  be a random sample from  $Unif(0, \theta)$ ,  $\theta > 0$  and the prior of  $\theta$  is  $\tau(\theta) = \theta \exp(-\theta)$ ,  $\theta > 0$ .  
 (a) Find the posterior distribution of  $\theta$  given  $X_1, \dots, X_n$ .  
 (b) Find the Bayes estimators of  $\theta$  under the loss functions  $L_1(\theta, a) = |\theta - a|$  and  $L_2(\theta, a) = (\theta - a)^2$ .  
**4.9.** We want to estimate the mean  $\theta$  of a Poisson distribution on the basis of a random sample  $X_1, \dots, X_n$  subject to the loss function  $L(\theta, a) = (\theta - a)^2$ ,  $\theta \in \Theta = (0, \infty)$  and  $a \in A = [0, \infty)$ . Assume that  $\theta$  has a  $\text{Gamma}(\alpha, \beta)$  prior (ie, it has the pdf)

$$\begin{aligned} \tau_{\alpha, \beta}(\theta) &= \{\Gamma(\alpha)\beta^\alpha\}^{-1} \theta^{\alpha-1} \exp(-\theta/\beta), \theta > 0, \text{ with} \\ E_{\alpha, \beta}(\theta) &= \alpha\beta, \quad E_{\alpha, \beta}(\theta^2) = \alpha(\alpha + 1)\beta^2. \end{aligned}$$

- (a) Find the posterior distribution of  $\theta$  given  $X_1, \dots, X_n$  and the Bayes rule with respect to  $\tau_{\alpha, \beta}$ .
- (b) Show that  $d_0(x) = \bar{X}_n$  is not a Bayes rule with respect to any prior on  $\Theta = (0, \infty)$ . However,  $d_0$  is a Bayes rule if  $\Theta = [0, \infty)$ .
- (c) Show that  $d_0$  is
  - (i) a limit of Bayes rules, (ii) a generalized Bayes rule with respect to  $\tau(\theta) = \theta^{-1}, \theta > 0$ , which is not a pdf, and (iii) an extended Bayes rule.
- 4.10.** Let  $X$  and  $Y$  be independent binomial  $\text{Bin}(n, \theta_1)$  and  $\text{Bin}(n, \theta_2)$  rv's. We want to estimate  $\theta_1 - \theta_2$  subject to the squared error loss  $L((\theta_1, \theta_2), a) = ((\theta_1 - \theta_2) - a)^2$ ,  $|a| \leq 1$ . Find the Bayes rule with respect to the prior with pdf  $\tau(\theta_1, \theta_2) = I_{(0,1)}(\theta_1)I_{(0,1)}(\theta_2)$  (ie,  $\theta_1$  and  $\theta_2$  have independent uniform distributions on  $(0, 1)$ ).
- 4.11.** Suppose that  $X$  has a pdf  $f(x, \theta)$ ,  $T$  is a sufficient statistics for  $\theta$ , and that the factorization theorem holds. Show that all Bayes rules are functions of  $T$ .
- 4.12.** Let  $X_1, \dots, X_n$  be a random sample from a Weibull distribution with pdf

$$f(x, \theta) = \theta a x^{\theta-1} \exp(-\theta x^\theta), \quad x > 0, \quad \theta > 0,$$

where  $a > 0$  is known. Find a sufficient statistic for  $\theta$ .

- 4.13.** Let  $X_1, \dots, X_n$  be a random sample from a distribution with pdf

$$f(x, \theta) = a(\theta)h(x)I(\theta_1 \leq x \leq \theta_2), \quad \theta = (\theta_1, \theta_2), \quad -\infty < \theta_1 < \theta_2 < \infty,$$

where  $h(x) > 0$  is a known function with  $\int_{-\infty}^{\infty} h(x) dx = 1$  and  $a(\theta) = 1 / \int_{\theta_1}^{\theta_2} h(x) dx$  for  $\theta = (\theta_1, \theta_2)$ . Find a two-dimensional sufficient statistics for  $\theta$  and apply your result to the special case of uniform distribution on  $[\theta_1, \theta_2]$ .

# Point Estimation in Parametric Models

## 5.1 Optimality Under Unbiasedness, Squared-Error Loss, UMVUE

In the general framework of statistical decision theory (see [Section 4.4](#)), the estimation problem is described by the triple  $(\{P_\theta, \theta \in \Theta\}, \mathcal{A}, L)$  where  $\mathcal{A} = \{g(\theta) : \theta \in \Theta\}$  and  $L(\theta, a)$  = loss due to estimating  $g(\theta)$  by  $a$ . For a decision rule  $d$ , which is an estimator  $T = T(X_1, \dots, X_n)$  based on the data  $\mathbf{X} = (X_1, \dots, X_n)$ , the risk is

$$R(\theta, T) = E_\theta[L(\theta, T(\mathbf{X}))] = \int_{\mathcal{X}} L(\theta, T(\mathbf{x}))f(\mathbf{x}; \theta) d\mathbf{x}$$

in the continuous case and analogously in the discrete case.

The concept of unbiasedness in estimation has been introduced in [Section 4.6.1](#). In a parametric family  $\{(\mathcal{X}, \mathcal{A}, P_\theta), \theta \in \Theta\}$ , a statistic  $T$  based on a random sample  $\mathbf{X}$  is an unbiased estimator of  $g(\theta)$  if  $E_\theta[T(\mathbf{X})] = g(\theta)$  for all  $\theta \in \Theta$ .

**Example 5.1.1.** Let  $(X_1, \dots, X_n)$  be a random sample from  $Bernoulli(\theta)$  and let  $g(\theta) = \theta(1 - \theta)$ . Let  $T = X_1(1 - X_2)$ . Then

$$E_\theta[T] = E_\theta[X_1(1 - X_2)] = \theta(1 - \theta) = g(\theta).$$

**Example 5.1.2.** Let  $(X_1, \dots, X_n)$  be a random sample from  $Unif(\theta)$  and let  $g(\theta) = \theta$ . Let  $T = 2X_1$ . Then

$$E_\theta[T] = E_\theta[2X_1] = 2E_\theta[X_1] = 2 \cdot \frac{\theta}{2} = \theta = g(\theta).$$

**Example 5.1.3.** Let  $(X_1, \dots, X_n)$  be a random sample from  $N(\mu, \sigma^2)$ , and let  $g_1(\mu, \sigma) = \mu$ , and  $g_2(\mu, \sigma) = \sigma^2$ . Let  $T_1 = \bar{X} = n^{-1} \sum_{i=1}^n X_i$  and  $T_2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$ . Then

$$\begin{aligned} E_{\mu, \sigma}[T_1] &= n^{-1} E_{\mu, \sigma} \left[ \sum_{i=1}^n X_i \right] = n^{-1} n \mu = \mu, \text{ and} \\ E_{\mu, \sigma}[T_2] &= (n-1)^{-1} E_{\mu, \sigma} \left[ \sum_{i=1}^n (X_i - \bar{X})^2 \right] \end{aligned}$$

$$\begin{aligned}
&= (n-1)^{-1} E_{\mu, \sigma} \left[ \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2 \right] \\
&= (n-1) \left[ n\sigma^2 - n \frac{\sigma^2}{n} \right] = \sigma^2.
\end{aligned}$$

Thus in these examples,  $X_1(1 - X_2)$  is an unbiased estimator of  $\theta(1 - \theta)$ ,  $2X_1$  is an unbiased estimator of  $\theta$ , and  $\bar{X}$  and  $(n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$  are unbiased estimators of  $\mu$  and  $\sigma^2$ , respectively.

**Theorem 5.1.1** (Rao-Blackwell Theorem). *Let  $A$  be a convex set in  $\mathbb{R}^k$ ,  $L(\theta, a)$  convex in  $a \in A$  for each  $\theta \in \Theta$  and  $T$  sufficient for  $\theta$  in  $\mathbf{X}$ . Then for any nonrandomized decision rule  $d$ , the nonrandomized decision rule  $d^*$  based on  $T$ , defined as*

$$d^*(t) = E[d(\mathbf{X})|T = t], \quad (1)$$

assuming that the expectation exists, is at least as good as  $d$  (ie,  $R(\theta, d^*)$  is either  $< R(\theta, d)$  or  $= R(\theta, d)$  for all  $\theta$ ).

*Proof.* Sufficiency of  $T$  makes  $d^*(T) = E[d(\mathbf{X})|T]$  free of  $\theta$  so that  $d^*(T)$  is a bona fide estimator. By Jensen's inequality (see [Section A.2](#)), for fixed  $t$ ,

$$E[L(\theta, d(\mathbf{X}))|T = t] \geq L(\theta, E[d(\mathbf{X})|T = t]) = L(\theta, d^*(t)) \quad \text{for all } \theta.$$

Hence, for each  $\theta$ ,

$$\begin{aligned}
R(\theta, d) &= E_\theta [L(\theta, d(\mathbf{X}))] = E_\theta E[L(\theta, d(\mathbf{X}))|T] \\
&\geq E_\theta [L(\theta, d^*(T))] = R(\theta, d^*)
\end{aligned}$$

showing that  $d^*$  is “at least as good” as  $d$ . □

*Remark 5.1.1.* If  $d$  is an unbiased estimator of  $g(\theta)$  (ie,  $E_\theta [d(\mathbf{X})] = g(\theta)$  for all  $\theta \in \Theta$ ), then

$$E_\theta [d^*(T)] = E_\theta [E(d(\mathbf{X})|T)] = E_\theta [d(\mathbf{X})] = g(\theta) \quad \text{for all } \theta.$$

Thus,  $d^*$  is also an unbiased estimator of  $g(\theta)$ .

*Remark 5.1.2.* If  $d$  is an unbiased estimator of  $g(\theta)$  then, for  $L(\theta, a) = \{a - g(\theta)\}^2$ ,

$$R(\theta, d) = E_\theta [\{d(\mathbf{X}) - g(\theta)\}^2] = \text{Var}_\theta [d(\mathbf{X})],$$

and  $R(\theta, d^*)$  is also  $\text{Var}_\theta [d^*(T)]$ . Hence, if  $d$  is an unbiased estimator of  $g(\theta)$ , then  $d^*$  is also an unbiased estimator of  $g(\theta)$ , the variance of which is uniformly bounded above by the variance of  $d$ .

The Rao-Blackwell Theorem is often stated in this form.

**Theorem 5.1.2** (Rao-Blackwell Theorem for Squared-Error Loss). *If  $d$  is an unbiased estimator of  $g(\theta)$  based on a sample  $\mathbf{X}$  from  $P_\theta$ ,  $\theta \in \Theta$ , and if  $T$  is sufficient for  $\theta$  in  $\mathbf{X}$ , then  $d^*$ , defined by*

$$d^*(t) = E[d(\mathbf{X})|T = t],$$

is also an unbiased estimator of  $g(\theta)$ , and

$$\text{Var}_\theta[d^*(T)] \leq \text{Var}_\theta[d(X)] \quad \text{for all } \theta \in \Theta.$$

(See [15, 16].)

Under squared-error loss, the Rao-Blackwell Theorem provides us with a method for “improving” upon an unbiased estimator, using a sufficient statistic, by constructing an unbiased estimator with uniformly smaller (or equal) variance. However, we do not know whether there is an unbiased estimator with even smaller variance than this “improved” estimator. More to the point, among many sufficient statistics, which one should we use in the Rao-Blackwell formula to find the best unbiased estimator? Such an estimator will be called the Uniformly Minimum Variance Unbiased Estimator (UMVUE).

Before searching for the UMVUE, we shall first show that it is unique and find a characterization for the UMVUE.

**Theorem 5.1.3.** *The UMVUE is unique.*

*Proof.* Suppose  $T_1$  and  $T_2$  are two distinct UMVUEs of  $g(\theta)$ , that is

$$\begin{aligned} E_\theta[T_1] &= E_\theta[T_2] = g(\theta), \text{ and} \\ \text{Var}_\theta[T_1] &= \text{Var}_\theta[T_2] := \sigma^2(\theta) \quad \text{for all } \theta. \end{aligned}$$

Then, by Cauchy-Schwartz,

$$\text{Cov}_\theta[T_1, T_2] \leq \{\text{Var}_\theta[T_1]\text{Var}_\theta[T_2]\}^{1/2} = \sigma^2(\theta).$$

Now let  $\bar{T} = (T_1 + T_2)/2$ . Then,  $\bar{T}$  is also an unbiased estimator of  $g(\theta)$  and

$$\text{Var}_\theta[\bar{T}] = \frac{1}{4}\{\text{Var}_\theta[T_1] + \text{Var}_\theta[T_2] + 2\text{Cov}_\theta[T_1, T_2]\} \leq \sigma^2(\theta) = \text{Var}_\theta[T_1].$$

If this inequality is strict, then the UMVUE property of  $T_1$  is violated. On the other hand, equality holds only if we have equality in the Cauchy-Schwartz, for which we need  $T_2 = a(\theta) + b(\theta)T_1$ , but then,

$$\sigma^2(\theta) = \text{Cov}_\theta[T_1, T_2] = \text{Cov}_\theta[T_1, a(\theta) + b(\theta)T_1] = b(\theta)\sigma^2(\theta),$$

so  $b(\theta) = 1$  and  $a(\theta) = 0$  since  $E_\theta[T_1] = E_\theta[T_2] = g(\theta)$ . Thus,  $T_1 = T_2$  showing that  $T_1$  is unique.  $\square$

**Definition 5.1.1.** Let  $\mathcal{U}$  denote the class of estimators  $U$  with  $E_\theta[U] = 0$  and  $E_\theta[U^2] < \infty$  for all  $\theta$ .

**Theorem 5.1.4.** *If  $T$  is an unbiased estimator of  $g(\theta)$ , then a necessary and sufficient condition for  $U$  to be the UMVUE of  $g(\theta)$  is  $\text{Cov}_\theta[T, U] = E_\theta[TU] = 0$  for all  $U \in \mathcal{U}$ .*

*Proof (Necessity).* Let  $T$  be UMVUE of  $g(\theta)$ . Then, for arbitrary  $U \in \mathcal{U}$  and  $\lambda \in \mathbb{R}$ ,  $T + \lambda U$  is also an unbiased estimator of  $g(\theta)$ . Hence,

$$\text{Var}_\theta[T] \leq \text{Var}_\theta[T + \lambda U] = \text{Var}_\theta[T] + \lambda^2\text{Var}_\theta[U] + 2\lambda\text{Cov}_\theta[T, U] \quad \text{for all } \lambda,$$

that is,  $f(\lambda) = \lambda^2 \text{Var}_\theta[U] + 2\lambda \text{Cov}_\theta[T, U] \geq 0$  for all  $\lambda$ . But if  $\text{Cov}_\theta[T, U] \neq 0$ , then

$$\min_{\lambda} f(\lambda) = -\frac{\{\text{Cov}_\theta[T, U]\}^2}{\text{Var}_\theta[U]} < 0.$$

To avoid contradiction, we therefore need  $\text{Cov}_\theta[T, U] = 0$ .

(Sufficiency). Suppose  $\text{Cov}_\theta[T, U] = \text{E}_\theta[TU] = 0$  for all  $U \in \mathcal{U}$ . We shall show that  $T$  is UMVUE of  $g(\theta)$ . Let  $T'$  be an arbitrary unbiased estimator of  $g(\theta)$  with  $\text{Var}_\theta[T'] < \infty$ . Then,  $T - T' \in \mathcal{U}$ , so that  $\text{E}_\theta[T(T - T')] = 0$  (ie,  $\text{E}_\theta[T^2] = \text{E}_\theta[TT']$ ). Since  $\text{E}_\theta[T] = g(\theta)$ , it follows that

$$\begin{aligned}\text{Var}_\theta[T] &= \text{E}_\theta[T^2] - \{g(\theta)\}^2 = \text{E}_\theta[TT'] - \{g(\theta)\}^2 \\ &= \text{Cov}_\theta[T, T'] \leq \{\text{Var}_\theta[T]\text{Var}_\theta[T']\}^{1/2},\end{aligned}$$

which implies  $\text{Var}_\theta[T] \leq \text{Var}_\theta[T']$ .

The main steps in finding the UMVUE of  $g(\theta)$  are

- (i) finding an unbiased estimator of  $g(\theta)$ , and
- (ii) using an *appropriate* sufficient statistic  $T$  to use in Eq. (1), the Rao-Blackwell formula  $d^*(T) = \text{E}[d(\mathbf{X})|T]$ .

The following example will illustrate these issues.

**Example 5.1.4.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $Bernoulli(\theta)$ ,  $0 < \theta < 1$ . We want to estimate  $g(\theta) = \theta^2$ .

It is easy to see that  $d(\mathbf{X}) = X_1 X_2$  is an unbiased estimator of  $\theta^2$ :

$$\text{E}_\theta[X_1 X_2] = \text{E}_\theta[X_1] \text{E}_\theta[X_2] = \theta^2 \quad \text{for all } \theta.$$

In many situations, as in this example, finding an unbiased estimator of  $g(\theta)$  is easy. Getting into the question of the choice of a sufficient statistic (Step (ii)), first note that there are many sufficient statistics of which we consider the following three:

$$\begin{aligned}T_1 &= (X_1 + X_2, X_3 + \dots + X_n), \quad T_2 = (X_1 + X_n, X_2 + \dots + X_{n-1}), \text{ and} \\ T_3 &= X_1 + \dots + X_n.\end{aligned}$$

It is easy to see that they are all sufficient statistics, which we leave as exercises. We now use the Rao-Blackwell formula, conditioning  $d(\mathbf{X}) = X_1 X_2$  by  $T_1, T_2, T_3$  given above:

- (a) We first note that

$$\begin{aligned}\text{E}[d(\mathbf{X})|T_1] &= \text{E}[X_1 X_2 | X_1 + X_2, X_3, \dots, X_n] \\ &= \text{E}[X_1 X_2 | X_1 + X_2] = X_1 X_2,\end{aligned}$$

because

$$\begin{aligned}\text{E}[X_1 X_2 | X_1 + X_2 = 0] &= 0 = X_1 X_2, \\ \text{E}[X_1 X_2 | X_1 + X_2 = 1] &= 0 = X_1 X_2, \text{ and} \\ \text{E}[X_1 X_2 | X_1 + X_2 = 2] &= 1 = X_1 X_2.\end{aligned}$$

**(b)** To find  $E[d(\mathbf{X})|T_2]$ , we first note that

$$\begin{aligned} E[X_1 X_2 | X_1 + X_n = r, X_2 + \dots + X_{n-1} = s] \\ = P[X_1 = X_2 = 1 | X_1 + X_n = r, X_2 + \dots + X_{n-1} = s] \\ := P_{rs}. \end{aligned}$$

Obviously,  $P_{rs} = 0$  if  $r = 0$  and/or  $s = 0$ . Next, for  $r = 1, 2$  and  $s = 1, \dots, n - 2$ ,

$$\begin{aligned} P_{rs} &= \frac{P[X_1 = X_2 = 1, X_n = r - 1, X_3 + \dots + X_{n-1} = s - 1]}{P[X_1 + X_n = r, X_2 + \dots + X_{n-1} = s]} \\ &= \frac{\theta^2 \theta^{r-1} (1-\theta)^{2-r} \binom{n-3}{s-1} \theta^{s-1} (1-\theta)^{n-s-2}}{\binom{2}{r} \theta^r (1-\theta)^{2-r} \binom{n-2}{s} \theta^s (1-\theta)^{n-s-2}} \\ &= \frac{rs}{2(n-2)}, \end{aligned}$$

after simplification, since  $\binom{2}{r} = 2/r$  for  $r = 1$  or  $r = 2$ . Thus,

$$E[d(\mathbf{X})|T_2] = \frac{(X_1 + X_n)(X_2 + \dots + X_{n-1})}{2(n-2)}.$$

**(c)** Note that

$$\begin{aligned} E[d(\mathbf{X})|T_3 = t] &= P[X_1 = X_2 = 1 | X_1 + \dots + X_n = t] \\ &= \frac{P[X_1 = X_2 = 1, X_3 + \dots + X_n = t - 2]}{P[X_1 + \dots + X_n = t]} \\ &= \frac{\theta^2 \binom{n-2}{t-2} \theta^{t-2} (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} \\ &= \frac{\binom{n-2}{t-2}}{\binom{n}{t}} = \frac{t(t-1)}{n(n-1)}. \end{aligned}$$

Thus,

$$E[d(\mathbf{X})|T_3] = \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n X_i - 1)}{n(n-1)}.$$

In applying the Rao-Blackwell method on  $d(\mathbf{X}) = X_1 X_2$ , using the sufficient statistics  $T_1$ ,  $T_2$ , and  $T_3$ , we have seen that  $E[d(\mathbf{X})|T_1] = d(\mathbf{X})$ , so the method provides an “improvement” in a trivial sense. On the other hand,  $E[d(\mathbf{X})|T_2]$  provides a real improvement, because  $d(\mathbf{X})$  is not a function of  $T_2$ , so the conditional distribution of  $d(\mathbf{X})$  given  $T_2$  is nondegenerate. Also, by direct calculation, we can verify that

$$\text{Var}_\theta E[d(\mathbf{X})|T_2] = \frac{\theta^2(1-\theta)}{2(n-2)}[1 + (n-1)\theta] < \theta^2(1-\theta)^2 = \text{Var}_\theta[d(\mathbf{X})]$$

for  $n \geq 3$ . Since  $d(\mathbf{X})$  is not a function of  $T_3$ , the estimator  $E[d(\mathbf{X})|T_3]$  also provides a real improvement over  $d(\mathbf{X})$ , although a direct calculation of the variance of  $E[d(\mathbf{X})|T_3]$  is somewhat messy.

Finally, we use the Rao-Blackwell method on the estimator  $E[d(\mathbf{X})|T_2]$ , using the sufficient statistic  $T_3$ . This leads to

$$\begin{aligned} & E\left[\frac{(X_1 + X_n)(X_2 + \cdots + X_{n-1})}{2(n-2)} \middle| X_1 + \cdots + X_n = t\right] \\ &= \sum_{r=1}^2 \frac{r(t-r)P[X_1 + X_n = r, X_2 + \cdots + X_{n-1} = t-r]}{2(n-2)\binom{n}{t}\theta^t(1-\theta)^{n-t}} \\ &= \sum_{r=1}^2 \frac{r(t-r)\binom{2}{r}\theta^r(1-\theta)^{2-r}\binom{n-2}{t-r}\theta^{t-r}(1-\theta)^{n-t+r-2}}{2(n-2)\binom{n}{t}\theta^t(1-\theta)^{n-t}} \\ &= \sum_{r=1}^2 \frac{r(t-r)\binom{2}{r}\binom{n-2}{t-r}}{2(n-2)\binom{n}{t}} = \frac{t(t-1)}{n(n-1)} \end{aligned}$$

since the numerator simplifies to  $2(n-2)\binom{n-2}{t-2}$ .

In summary, we have seen that starting with an unbiased estimator  $d(\mathbf{X}) = X_1 X_2$  of  $g(\theta) = \theta^2$  and using the Rao-Blackwell method with sufficient statistics  $T_1 = (X_1 + X_2, X_3 + \cdots + X_n)$ ,  $T_2 = (X_1 + X_n, X_2 + \cdots + X_{n-1})$ , and  $T_3 = X_1 + \cdots + X_n$ , we obtained:

- (a)  $E[d(\mathbf{X})|T_1] = d(\mathbf{X})$ , which is an “improvement” in a trivial sense,
- (b)  $E[d(\mathbf{X})|T_2] = \frac{(X_1 + X_n)(X_2 + \cdots + X_{n-1})}{2(n-2)}$ , which is a real improvement,
- (c)  $E[d(\mathbf{X})|T_3] = \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n X_i - 1)}{n(n-1)}$ , which is also an improvement, and moreover,
- (d)  $E[E[d(\mathbf{X})|T_2]|T_3] = E[d(\mathbf{X})|T_3]$ .

Thus,  $E[d(\mathbf{X})|T_2]$  and  $E[d(\mathbf{X})|T_3]$  are successive improvements on  $d(\mathbf{X})$ .

The question still remains: can  $E[d(\mathbf{X})|T_3]$  be further improved? The key property of a sufficient statistic addressing this question is the property of completeness, as defined below.

**Definition 5.1.2.** A sufficient statistic  $T$  for a parameter  $\theta \in \Theta$  is said to be complete (or boundedly complete) if for every real-valued (bounded real-valued) function  $\varphi$ ,  $E_\theta[\varphi(T)] = 0$  for all  $\theta$  implies  $\varphi(t) = 0$  for all  $t \notin N$  where  $P_\theta[T \in N] = 0$ .

Completeness requires the family of distributions  $\{P_\theta^T\}$  of  $T$  to be sufficiently rich, so that the condition  $E_\theta[\varphi(T)] = 0$  for all  $\theta$  forces  $\varphi(t)$  to be identically zero for all practical purposes. One can think of it as a condition on the pdf's  $\{f_T(\cdot, \theta), \theta \in \Theta\}$  of  $T$  having full rank, so that

$$E_\theta[\varphi(T)] = \int \varphi(t)f_T(t, \theta) dt = 0 \quad \text{for all } \theta,$$

that is,  $\varphi(\cdot)$  is “orthogonal” to  $f_T(\cdot, \theta)$  for all  $\theta \in \Theta \Rightarrow \varphi(\cdot)$  must be a zero function with probability 1.

How to check that a sufficient statistic is complete?

1. If a convergent power series  $\sum_n a_n z^n = 0$  for all  $z$  in some open interval, then each coefficient  $a_n$  must be zero. This fact can be used to prove completeness of some

important families of discrete distributions such as  $\{Bin(n, p), n \text{ fixed}, 0 < p < 1\}$ ,  $\{Poi(\theta), \theta > 0\}$ , etc.

2. If  $X_1, \dots, X_n$  are iid as  $Unif(0, \theta)$ , then  $T = \max\{X_1, \dots, X_n\}$  is sufficient for  $\theta$  in  $\mathbf{X} = (X_1, \dots, X_n)$ . The pdf of  $T$  corresponding to  $\theta$  is

$$f_T(t, \theta) = n\theta^{-n}t^{n-1}I_{(0,\theta)}(t).$$

Hence,  $E_\theta[\varphi(T)] = 0$  for all  $\theta$  implies  $\int_0^\theta \varphi(t)t^{n-1} dt = 0$  for all  $\theta$ . If  $\varphi$  is continuous, then  $\varphi(t) = 0$  for all  $t$  by the Fundamental Theorem of Calculus, but even without continuity, the result holds with probability 1, of which the proof needs more advanced analysis. This shows that  $T = \max\{X_1, \dots, X_n\}$  is a complete sufficient statistic for  $\theta$  in  $\mathbf{X}$ .

3. If  $X_1, \dots, X_n$  is a random sample from a regular  $k$ -parameter exponential family:

$$f(x, \theta) = \exp \left[ \sum_{j=1}^k \theta_j T_j(\mathbf{x}) + S(\mathbf{x}) + d(\theta) \right] I_A(\mathbf{x}),$$

where  $A$  is the support of the distribution, then

$$T_j = \sum_{i=1}^n T_j(x_i), \quad j = 1, \dots, k$$

are jointly sufficient for  $\theta$  in  $\mathbf{X} = (X_1, \dots, X_n)$  and the joint distribution of  $\mathbf{T} = (T_1, \dots, T_k)$  also belongs to a regular  $k$ -parameter exponential family. If  $\Theta = \{\theta \in \mathbb{R}^k : \int_X f(\mathbf{x}, \theta) d\mathbf{x} < \infty\}$  contains a  $k$ -dimensional open rectangle, then  $\mathbf{T}$  is a complete sufficient statistic for  $\theta$  in  $\mathbf{X}$ . This takes care of many important situations.

With the concept of completeness, we now know how to choose the *appropriate sufficient statistic* in the Rao-Blackwell formula (1) that will lead to the UMVUE.

**Theorem 5.1.5** (Lehmann-Scheffé). *If  $d(\mathbf{X})$  is an unbiased estimator of  $g(\theta)$  and  $T$  is a complete sufficient statistic for  $\theta$  in  $\mathbf{X}$ , then*

$$d^*(T) = E[d(\mathbf{X})|T]$$

is the (essentially unique) UMVUE of  $g(\theta)$ .

*Proof.* For any other unbiased estimator  $d_1(\mathbf{X})$  of  $g(\theta)$ , consider the unbiased estimator  $d_1^*(T) = E[d_1(\mathbf{X})|T]$  which must have variance  $\leq$  that of  $d_1$  by the Rao-Blackwell Theorem, but by completeness,  $d_1^*(T) = d^*(T)$  with probability 1. Therefore,

$$\text{Var}_\theta[d_1(\mathbf{X})] \geq \text{Var}_\theta[d_1^*(T)] = \text{Var}_\theta[d^*(T)]. \quad \square$$

In the above example,  $T_3 = \sum_{i=1}^n X_i$  is a complete sufficient statistic. Hence,

$$E[d(\mathbf{X})|T_3] = \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n X_i - 1)}{n(n-1)}$$

is the UMVUE of  $g(\theta) = \theta^2$ . Also, note that the sufficient statistic  $T_2 = (X_1 + X_n, X_2 + \dots + X_{n-1})$  is *not complete*, because the expectation of

$$\varphi(T_2) = \frac{1}{2}(X_1 + X_n) - \frac{1}{n-2}(X_2 + \dots + X_{n-1})$$

is 0 for all  $\theta \in (0, 1)$ . We end this section with a very interesting property of completeness.

### Ancillarity and Completeness

**Definition 5.1.3.** A statistic whose distribution does not depend on  $\theta$  is called an ancillary statistic.

**Theorem 5.1.6** (Basu's Theorem [17]). *Suppose that the distribution of  $X$  belongs to the family  $\{P_\theta, \theta \in \Theta\}$ . If  $T = T(X)$  is a complete sufficient statistic for  $\theta$  and if  $V = V(X)$  is an ancillary statistic, then  $V$  is independent of  $T$ .*

*Proof.* The probability  $p_A = P_\theta[V \in A]$  is independent of  $\theta$  for all  $A$  by ancillarity of  $V$ . Let  $\eta_A(t) = P[V \in A|T = t]$ , which is also independent of  $\theta$  since  $T$  is sufficient, and  $E_\theta[\eta_A(T)] = E_\theta P[V \in A|T] = P_\theta[V \in A] = p_A$ . Thus

$$E_\theta[\eta_A(T) - p_A] = 0 \quad \text{for all } \theta.$$

Since  $T$  is complete, this implies that for all  $A$ ,  $\eta_A(T) = p_A$ , that is

$$P[V \in A|T] = P[V \in A] \text{ w.p. 1.}$$

Hence,  $V$  is independent of  $T$ . □

## 5.2 Lower Bound for the Variance of an Unbiased Estimator

### 5.2.1 The Information Inequality: Cramèr-Rao Lower Bound

For brevity of notation, we shall write

$$\begin{aligned} \frac{\partial f(x, \theta)}{\partial \theta} &= \dot{f}(x, \theta), & \frac{\partial^2 f(x, \theta)}{\partial \theta^2} &= \ddot{f}(x, \theta), \\ \log f(x, \theta) &= \ell(x, \theta), & \frac{\partial \ell(x, \theta)}{\partial \theta} &= \dot{\ell}(x, \theta), & \frac{\partial^2 \ell(x, \theta)}{\partial \theta^2} &= \ddot{\ell}(x, \theta). \end{aligned}$$

These notations will also be used in all subsequent discussions.

Regularity conditions on  $\{f(x, \theta), \theta \in \Theta \subset \mathbb{R}\}$ :

1. The parameter space  $\Theta$  is an open interval and the set  $S = \{x: f(x, \theta) > 0\}$  does not depend upon  $\theta$ .
2. For all  $x \in S$  and  $\theta \in \Theta$ ,  $\dot{\ell}(x, \theta)$  exists and is finite.
3. For any statistic  $T$  such that  $E_\theta[|T|] < \infty$  for all  $\theta \in \Theta$ ,

$$\begin{aligned} \frac{d}{d\theta} \int T(x)f(x, \theta) dx \Big|_{\theta=\theta_0} &= \int T(x)\dot{f}(x, \theta_0) dx \\ &= \int T(x)\dot{\ell}(x, \theta_0)f(x, \theta_0) dx \end{aligned}$$

whenever the right-hand side is finite. In other words,

$$\frac{d}{d\theta} \mathbb{E}_\theta[T(X)] \Big|_{\theta = \theta_0} = \mathbb{E}_{\theta_0}[T(X)\dot{\ell}(X, \theta_0)].$$

These conditions are satisfied in a regular exponential family.

*Remark 5.2.1.* Write  $\mathbb{E}_\theta[T(X)] = g(\theta) + b(\theta)$ , where  $b(\theta)$  is the bias of  $T$  at  $\theta$  as an estimator of  $g(\theta)$ . Then, Condition 3 becomes

$$g'(\theta) + b'(\theta) = \mathbb{E}_\theta[T(X)\dot{\ell}(X, \theta)].$$

*Remark 5.2.2.*

(a) Taking  $T(x) \equiv 1$  in Condition 3, we have  $\mathbb{E}_\theta[\dot{\ell}(X, \theta)] = 0$ .

(b) Sometimes (as in regular exponential families), a stronger version of Condition 3,

$$\frac{d^2}{d\theta^2} \int T(x)f(x, \theta) dx \Big|_{\theta = \theta_0} = \int T(x)\ddot{f}(x, \theta_0) dx$$

holds for  $T(x) \equiv 1$ , so that  $\int \ddot{f}(x, \theta) dx = 0$ . Since

$$\ddot{f}(x, \theta) = [\ddot{\ell}(x, \theta) + \{\dot{\ell}(x, \theta)\}^2]f(x, \theta),$$

this implies that

$$\text{Var}_\theta[\dot{\ell}(X, \theta)] = \mathbb{E}_\theta[\{\dot{\ell}(X, \theta)\}^2] = -\mathbb{E}_\theta[\ddot{\ell}(X, \theta)].$$

**Definition 5.2.1** (Fisher-Information). Condition 2 allows us to define

$$I(\theta_0) = \mathbb{E}_{\theta_0}[\{\dot{\ell}(X, \theta_0)\}^2], \tag{2}$$

which is called the Fisher-Information in  $X$  for the family  $\{f(x, \theta), \theta \in \Theta\}$  at  $\theta_0$ .

*Remark 5.2.3.* Note that  $0 \leq I(\theta) \leq \infty$ . By Condition 3,

$$I(\theta) = \text{Var}_\theta[\dot{\ell}(X, \theta)],$$

and the stronger version of Condition 3 implies

$$I(\theta) = \mathbb{E}_\theta[\{\dot{\ell}(X, \theta)\}^2] = \text{Var}_\theta[\dot{\ell}(X, \theta)] = \mathbb{E}_\theta[-\ddot{\ell}(X, \theta)].$$

**Theorem 5.2.1** (Cramér-Rao Inequality [15, 18]). *Suppose  $T = T(X)$  has  $\mathbb{E}_\theta[T] = g(\theta) + b(\theta)$  and  $\text{Var}_\theta[T] < \infty$ . Then, under the regularity Conditions 1, 2, and 3, we have*

$$\text{Var}_\theta[T] \geq \frac{\{g'(\theta) + b'(\theta)\}^2}{I(\theta)}.$$

The right-hand side of this inequality is known as the *Cramér-Rao lower bound*.

*Proof.* Fix  $\theta \in \Theta$  and let  $S = S(X, \theta) = \dot{\ell}(X, \theta)$ . By Remark 5.2.2,  $\mathbb{E}_\theta[S] = 0$  and  $\text{Var}_\theta[S] = I(\theta)$ . Hence,

$$\text{Cov}_\theta[T, S] = \mathbb{E}_\theta[TS] = \mathbb{E}_\theta[T(X)\dot{\ell}(X, \theta)] = g'(\theta) + b'(\theta),$$

as shown earlier. But  $\text{Cov}_\theta^2[T, S] \leq \text{Var}_\theta[T]\text{Var}_\theta[S]$ . Thus,

$$\text{Var}_\theta[T] \geq \frac{\text{Cov}_\theta^2[T, S]}{\text{Var}_\theta[S]} = \frac{\{g'(\theta) + b'(\theta)\}^2}{I(\theta)}. \quad \square$$

*Remark 5.2.4.*

- (a) If  $\mathbf{X} = (X_1, \dots, X_n)$  is a random sample of size  $n$  from a population with pdf (or pmf)  $f(x, \theta)$ , then  $\dot{\ell}(\mathbf{X}, \theta) = \sum_{i=1}^n \dot{\ell}(X_i, \theta)$ . Hence, the Fisher-information in  $\mathbf{X} = (X_1, \dots, X_n)$  for the family  $\{f(\mathbf{x}, \theta): \theta \in \Theta\}$  is

$$\text{Var}_\theta \left[ \sum_{i=1}^n \dot{\ell}(X_i, \theta) \right] = n\text{Var}_\theta[\dot{\ell}(X_1, \theta)] = nI(\theta).$$

Thus the information inequality for  $\text{Var}_\theta[T]$  of an estimator  $T$  based on a random sample of size  $n$  from  $f(x, \theta)$  is

$$\text{Var}_\theta[T] \geq \frac{\{g'(\theta) + b'(\theta)\}^2}{nI(\theta)}.$$

- (b) If  $T$  is an unbiased estimator of  $\theta$  based on a random sample of size  $n$ , then

$$\text{Var}_\theta[T] \geq \frac{1}{nI(\theta)}.$$

### 5.2.2 Effect of Reparametrization on Fisher-Information in Exponential Families

Consider the exponential family of pdf's (pmf's):

$$f(x, \theta) = \exp[c(\theta)T(x) + d(\theta) + h(x)]I_A(x),$$

where  $c$  is one-to-one and twice differentiable. If  $f(x, \theta)$  is reparametrized by  $\eta = c(\theta)$  and rewritten as

$$g(x, \eta) = \exp[\eta T(x) + d_0(\eta) + h(x)]I_A(x)$$

with  $d_0(\eta) = d[c^{-1}(\eta)]$ , then the Fisher-Information for the family  $\{g(x, \eta): \eta \in c(\Theta)\}$  is  $I_g(\eta) = -d_0''(\eta)$ . How does this  $I_g$  relate to the Fisher-Information  $I_f(\theta) = E_\theta[-\ddot{\ell}(X, \theta)]$  for the original family? The answer is:

$$I_f(\theta) = (d_0''(\eta)|_{\eta=c(\theta)})\{c'(\theta)\}^2,$$

the proof of which is left as an exercise.

### 5.2.3 Information Inequality in Multiparameter Families

We now consider the multiparameter case in which we observe a random variable  $X$  following a distribution with pdf (or pmf) belonging to a family  $\{f(x, \theta): \theta \in \Theta\}$ , where  $\theta = (\theta_1, \dots, \theta_k)$  and  $\Theta \subset \mathbb{R}^k$ . Our aim is to estimate a function  $g(\theta)$  of the parameter vector  $\theta$  such as  $g(\theta) = \theta_r$ , or, as in the pdf of  $N(\theta_1, \theta_2)$  with  $\theta_1 \neq 0$ ,  $g(\theta) = \sqrt{\theta_2}/\theta_1$ , the coefficient of variation. For a statistic  $T = T(X)$ , estimating  $g(\theta)$  with a bias  $b(\theta)$  (ie,  $E_\theta[T(X)] = g(\theta) + b(\theta)$ ), we now look for a lower bound of  $\text{Var}_\theta[T]$ . The notations introduced earlier are now extended as:

$$\begin{aligned}\frac{\partial f(x, \theta)}{\partial \theta_r} &= \dot{f}_r(x, \theta), \quad \frac{\partial^2 f(x, \theta)}{\partial \theta_r \partial \theta_s} = \ddot{f}_{rs}(x, \theta), \\ \log f(x, \theta) = \ell(x, \theta) \text{ as before}, \quad \frac{\partial \ell(x, \theta)}{\partial \theta_r} &= \dot{\ell}_r(x, \theta), \quad \frac{\partial^2 \ell(x, \theta)}{\partial \theta_r \partial \theta_s} = \ddot{\ell}_{rs}(x, \theta)\end{aligned}$$

for  $r, s = 1, \dots, k$ .

The Cramèr-Rao Information Inequality is simply a restatement of the inequality  $\rho^2(U, V) \leq 1$ , where  $\rho(U, V)$  is the correlation between  $U$  and  $V$  with  $U = T(X)$  and  $V = \dot{\ell}(X, \theta)$ .

In the multiparameter case, we extend the concept of correlation between  $U$  and  $V$  to the

“correlation” between  $U = T(X)$  and  $V^\top = (\dot{\ell}_1(X, \theta), \dots, \dot{\ell}_k(X, \theta))$

which we shall call the multiple correlation of  $U$  on  $V$ .

Let  $V^\top = (V_1, \dots, V_k)$  be a  $k$ -dim random vector with covariance matrix

$$\Sigma = \text{Cov}[V, V] = E\left[\{V - E(V)\}\{V - E(V)\}^\top\right],$$

and let  $U$  be a random variable with finite second moment,

$$\gamma_i = \text{Cov}[U, V_i], \quad i = 1, \dots, k, \text{ and } \gamma^\top = (\gamma_1, \dots, \gamma_k).$$

**Definition 5.2.2.** If  $\Sigma = \text{Cov}[V, V]$  is positive definite, then the multiple correlation of  $U$  on  $V$  is defined as:

$$\rho^* = \rho_{U \cdot V}^* = \max_{a_1, \dots, a_k} \text{Corr}\left[U, \sum_{i=1}^k a_i V_i\right] = \max_{a_1, \dots, a_k} \frac{\mathbf{a}^\top \boldsymbol{\gamma}}{\{\text{Var}[U] \cdot \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{a}\}^{1/2}}.$$

Note that

- (i) when  $\Sigma$  is singular and is of rank  $r$ , we can always find a  $r$ -dim random vector  $V^*$ , which is a linear function of  $V$ , whose covariance matrix is positive definite, and
- (ii) we can assume without loss of generality that  $\mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{a} = 1$ .

*Justification of (i).* By spectral decomposition,  $\boldsymbol{\Sigma} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top$  where  $\mathbf{Q}$  is an orthogonal matrix with columns  $\mathbf{q}_1, \dots, \mathbf{q}_k$  (the eigenvectors of  $\boldsymbol{\Sigma}$ ) and  $\boldsymbol{\Lambda}$  is a diagonal matrix whose diagonal elements  $\lambda_1 \geq \dots \geq \lambda_k$  are the eigenvalues of  $\boldsymbol{\Sigma}$ . If  $\boldsymbol{\Sigma}$  is of rank  $r < k$ , then only the first

$r$  diagonal elements of  $\Sigma$  are nonzero and the rest are zeros. Let  $V^* = [\mathbf{q}_1^\top V, \dots, \mathbf{q}_r^\top V]^\top = \mathbf{Q}_0^\top V$ , where  $\mathbf{Q}_0 = [\mathbf{q}_1, \dots, \mathbf{q}_r]$  is a  $k \times r$  matrix. Then,

$$\text{Cov}[V_i^*, V_j^*] = \text{Cov}[\mathbf{q}_i^\top V, \mathbf{q}_j^\top V] = \mathbf{q}_i^\top \Sigma \mathbf{q}_j = \begin{cases} \lambda_i & i=j \\ 0 & i \neq j \end{cases}$$

Since  $\Sigma$  is of rank  $r$ ,  $\lambda_1 \geq \dots \geq \lambda_r > 0$  and  $\lambda_{r+1} = \dots = \lambda_k = 0$ . Consequently, the covariance matrix  $\Sigma^*$  of  $V^* = \mathbf{Q}_0^\top V$  is a diagonal matrix with diagonal elements  $\lambda_1, \dots, \lambda_k$  (ie,  $\Sigma^*$  is positive definite).

*Justification of (ii).* Correlation coefficient is invariant under scale change. Hence, for any vector  $\mathbf{a} \in \mathbb{R}^k$ , if we take  $\mathbf{a}_0 = \mathbf{a}/\sqrt{\mathbf{a}^\top \Sigma \mathbf{a}}$ , then

$$\text{Var}[\mathbf{a}_0^\top V] = 1 \text{ and } \text{Corr}[U, \mathbf{a}^\top V] = \text{Corr}[U, \mathbf{a}_0^\top V].$$

**Proposition 5.2.1.** Let  $U$  and  $(V_1, \dots, V_k)$  have finite second moments and let  $\gamma_i = \text{Cov}[U, V_i]$ ,  $i = 1, \dots, k$ , and  $\Sigma = \text{Cov}[V, V]$ . Then,

$$\rho^* = \frac{\boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma}}{\text{Var}[U]}.$$

*Proof.* Since  $\rho^* = \max_{a_1, \dots, a_k} \frac{\mathbf{a}^\top \boldsymbol{\gamma}}{\text{Var}[U]}$  and  $\mathbf{a}^\top \Sigma \mathbf{a} = 1$ , we need to maximize  $\mathbf{a}^\top \boldsymbol{\gamma}$  subject to the condition  $\mathbf{a}^\top \Sigma \mathbf{a} = 1$ . Using Lagrange's undetermined multiplier, we maximize

$$\mathbf{a}^\top \boldsymbol{\gamma} - \frac{1}{2} \lambda \mathbf{a}^\top \Sigma \mathbf{a}$$

and then find  $\lambda$  using  $\mathbf{a}^\top \Sigma \mathbf{a} = 1$ . To this end,

$$\begin{aligned} 0 &= \frac{\partial}{\partial a_r} \left[ \sum_{r=1}^k a_r \gamma_r - \frac{1}{2} \lambda \sum_{r=1}^k \sum_{s=1}^k a_r a_s \sigma_{rs} \right] \\ &= \gamma_r - \frac{1}{2} \lambda \left[ 2a_r \sigma_{rr} + 2 \sum_{r \neq s=1}^k a_s \sigma_{rs} \right] \\ &= \gamma_r - \lambda \sum_{s=1}^k a_s \sigma_{rs}, \quad r = 1, \dots, k \end{aligned}$$

implies  $\Sigma \mathbf{a} = \lambda^{-1} \boldsymbol{\gamma}$  (ie,  $\mathbf{a} = \lambda^{-1} \Sigma^{-1} \boldsymbol{\gamma}$ ).

To find  $\lambda$ , we now have the equation

$$1 = \mathbf{a}^\top \Sigma \mathbf{a} = \lambda^{-2} (\boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma}) = \lambda^{-2} \boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma},$$

with solutions:  $\lambda = \pm(\boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma})^{1/2}$ . Hence,

$$\mathbf{a}^* = \lambda^{-1} \Sigma^{-1} \boldsymbol{\gamma} = \frac{\pm \boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma}}{\sqrt{\boldsymbol{\gamma}^\top \Sigma^{-1} \boldsymbol{\gamma}}} \text{ and}$$

$$\mathbf{a}^T \boldsymbol{\gamma} = \pm \frac{\boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma}}{\sqrt{\boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma}}} = \pm \sqrt{\boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma}}$$

are the maximizers of  $\mathbf{a}^T \boldsymbol{\gamma}$  subject to  $\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} = 1$ . Hence,

$$\begin{aligned}\rho^{*2} &= \max_{\{a_1, \dots, a_k\}, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} = 1} \text{Corr}^2 \left[ U, \sum_{i=1}^k a_i V_i \right] = \max_{a_1, \dots, a_k} \left( \frac{\mathbf{a}^T \boldsymbol{\gamma}}{\sqrt{\text{Var}[U]}} \right)^2 \\ &= \frac{(\mathbf{a}^T \boldsymbol{\gamma})^2}{\text{Var}[U]} = \frac{\boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma}}{\text{Var}[U]}. \quad \square\end{aligned}$$

Getting back to  $U = T(X)$  and  $V^T = (\dot{\ell}_1(X, \theta), \dots, \dot{\ell}_k(X, \theta))$ , and recognizing that the multiple correlation  $\rho_{U,V}^*$ , like a correlation lies between  $-1$  and  $+1$ , so that

$$\rho^{*2} = \frac{\boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma}}{\text{Var}[U]} \leq 1, \quad \text{ie, } \text{Var}[U] \geq \boldsymbol{\gamma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\gamma},$$

we arrive at a lower bound for  $\text{Var}[U] = \text{Var}[T(X)]$ . To make this lower bound explicit, we need to find  $\boldsymbol{\gamma}^T = (\gamma_1, \dots, \gamma_k)$ , where  $\gamma_j = \text{Cov}_{\theta} [T(X), \dot{\ell}_j(X, \theta)]$ , and  $\boldsymbol{\Sigma} = ((\text{Cov}_{\theta} [\dot{\ell}_r(X, \theta), \dot{\ell}_s(X, \theta)]))$ .

As in the single-parameter case, here also we assume that the following regularity conditions hold for the family  $\{f(x, \theta), \theta \in \Theta \subset \mathbb{R}^k\}$ :

1. The parameter space  $\Theta$  is an open interval and the set  $S = \{x: f(x, \theta) > 0\}$  does not depend upon  $\theta$ .
2. For all  $x \in S$  and  $\theta \in \Theta$ ,  $\dot{\ell}_r(x, \theta)$  exists and is finite for  $r = 1, \dots, k$ .
3. For any statistic  $T$  such that  $E_{\theta}[|T|] < \infty$  for all  $\theta \in \Theta$ ,

$$\begin{aligned}\frac{\partial}{\partial \theta_r} \int T(x)f(x, \theta) dx \Big|_{\theta=\theta_0} &= \int T(x)\dot{f}_r(x, \theta_0) dx \\ &= \int T(x)\dot{\ell}_r(x, \theta_0)f(x, \theta_0) dx\end{aligned}$$

for  $r = 1, \dots, k$  whenever the right-hand side is finite. In other words,

$$\frac{\partial}{\partial \theta_r} E_{\theta}[T(X)] \Big|_{\theta=\theta_0} = E_{\theta_0}[T(X)\dot{\ell}_r(X, \theta_0)], \quad r = 1, \dots, k.$$

As in the single-parameter case, these conditions are also satisfied in a regular  $k$ -parameter exponential family.

Condition 3 becomes

$$E_{\theta}[T(X)\dot{\ell}_r(X, \theta)] = \frac{\partial}{\partial \theta_r}[g(\theta) + b(\theta)],$$

when  $T(X)$  is an estimator of  $g(\theta)$  with bias  $b(\theta)$ . Taking  $T(x) \equiv 1$  in Condition 3, here we have  $E_\theta[\dot{\ell}_r(X, \theta)] = 0$ ,  $r = 1, \dots, k$ , and a stronger version of Condition 3 (which holds in exponential families), namely,

$$\frac{\partial^2}{\partial \theta_r \partial \theta_s} \int T(x)f(x, \theta) dx \Big|_{\theta=\theta_0} = \int T(x)\ddot{f}_{rs}(x, \theta) dx$$

holds for  $T(x) \equiv 1$ , so that  $\int \ddot{f}_{rs}(x, \theta) dx = 0$ . Since

$$\ddot{f}_{rs}(x, \theta) = [\ddot{\ell}_{rs}(x, \theta) + \dot{\ell}_r(x, \theta)\dot{\ell}_s(x, \theta)]f(x, \theta),$$

this implies

$$\text{Cov}_\theta[\dot{\ell}_r(X, \theta), \dot{\ell}_s(X, \theta)] = E_\theta[\dot{\ell}_r(X, \theta)\dot{\ell}_s(X, \theta)] = -E_\theta[\ddot{\ell}_{rs}(X, \theta)].$$

**Definition 5.2.3** (Information Matrix). By Condition 2,  $\dot{\ell}_r(x, \theta)$  exist and is finite for  $r = 1, \dots, k$ , so we can define the  $k \times k$  matrix  $I(\theta_0) = ((I_{rs}(\theta_0)))$  where  $I_{rs}(\theta_0) = E_{\theta_0}[\dot{\ell}_r(X, \theta_0)\dot{\ell}_s(X, \theta_0)]$ , which is called the Information Matrix in  $X$  for the family  $\{f(x, \theta), \theta \in \Theta \subset \mathbb{R}^k\}$  at  $\theta_0$ .

By Condition 3,  $I_{rs}(\theta) = \text{Cov}_\theta[\dot{\ell}_r(X, \theta), \dot{\ell}_s(X, \theta)]$  and the stronger version of Condition 3 implies  $I_{rs}(\theta) = E_\theta[-\ddot{\ell}_{rs}(X, \theta)]$  as shown above.

We now go back to the inequality  $\text{Var}[T(X)] \geq \gamma^\top \Sigma^{-1} \gamma$ , where

$$\begin{aligned} \gamma^\top &= (\text{Cov}_\theta[T(X), \dot{\ell}_1(X, \theta)], \dots, \text{Cov}_\theta[T(X), \dot{\ell}_k(X, \theta)]) \\ &= \left( \frac{\partial}{\partial \theta_1} E_\theta[T(X)], \dots, \frac{\partial}{\partial \theta_r} E_\theta[T(X)] \right) \\ &= \nabla^\top [g(\theta) + b(\theta)]. \end{aligned}$$

by Condition 3. We thus arrive at the following theorem.

**Theorem 5.2.2** (Information Lower Bound in Multiparameter Case). *Suppose  $T = T(X)$  is an estimator of  $g(\theta)$  with bias  $b(\theta)$  and finite variance, that is,  $E_\theta[T] = g(\theta) + b(\theta)$  and  $\text{Var}_\theta[T] < \infty$  where the family  $\{f(x, \theta), \theta \in \Theta \subset \mathbb{R}^k\}$  satisfies regularity Conditions 1, 2, and 3. Then,*

$$\text{Var}_\theta[T] \geq \{ \nabla^\top [g(\theta) + b(\theta)] \} I(\theta)^{-1} \{ \nabla [g(\theta) + b(\theta)] \}.$$

**Example 5.2.1.** In a random sample  $(X_1, \dots, X_n)$  from  $N(\mu, \sigma^2)$ , let

$$\begin{aligned} \bar{X} &= n^{-1} \sum_{i=1}^n X_i, \quad W = n^{-1} \sum_{i=1}^n (X_i - \mu)^2, \text{ and} \\ S^2 &= (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2. \end{aligned}$$

We want UMVUEs of  $\sigma^2$  and  $\sigma$  when  $\mu$  is known and when  $\mu$  is unknown.

*Solution.* When  $\mu$  is known,  $W$  is a complete sufficient statistic for  $\sigma^2$ , and when  $\mu$  is unknown,  $(\bar{X}, S^2)$  is a complete sufficient statistic for  $(\mu, \sigma^2)$  in  $(X_1, \dots, X_n)$ . Therefore, it is enough to find unbiased estimators of  $\sigma^2$  and  $\sigma$ , which are functions of  $W$  when  $\mu$  is known and functions of  $(\bar{X}, S^2)$  when  $\mu$  is unknown, because these estimators would then be UMVUEs by the Lehmann-Scheffé Theorem. We now obtain such unbiased estimators.

When  $\mu$  is known,

- (a)  $W$  is an unbiased estimator of  $\sigma^2$ , since  $E_\sigma[(X_i - \mu)^2] = \sigma^2$ , and
- (b) starting with  $\sqrt{W}$  as a natural estimator of  $\sigma$ , we find

$$E_\sigma[\sqrt{W}] = \sigma n^{-1/2} E[(\chi_n^2)^{1/2}] = \frac{\sigma}{\sqrt{n}} \frac{\sqrt{2}\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})},$$

so that

$$T_0 = c_n \sqrt{W}, \text{ with } c_n = \sqrt{\frac{n}{2}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n+1}{2})},$$

is the UMVUE of  $\sigma$ .

When  $\mu$  is unknown,

- (a) recall from [Proposition 2.2.5](#) that

$$\frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2,$$

so that  $E_\sigma[S^2] = \frac{\sigma^2}{n-1} E[\chi_{n-1}^2] = \sigma^2$ , making  $S^2$  the UMVUE of  $\sigma^2$ , and

- (b) starting with  $S$  as an estimator of  $\sigma$ , we see that

$$E_\sigma[S] = \frac{\sigma}{\sqrt{n-1}} E[(\chi_{n-1}^2)^{1/2}] = \frac{\sigma}{\sqrt{n-1}} \frac{\sqrt{2}\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})},$$

so that

$$T_1 = c_{n-1} S, \text{ with } c_{n-1} = \sqrt{\frac{n-1}{2}} \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n}{2})},$$

is the UMVUE of  $\sigma$ .

In summary, we have

1.  $W$  and  $T_0 = c_n \sqrt{W}$  are the UMVUEs of  $\sigma^2$  and  $\sigma$ , respectively, when  $\mu$  is known, and
2.  $S^2$  and  $T_1 = c_{n-1} S$  are the UMVUEs of  $\sigma^2$  and  $\sigma$ , respectively, when  $\mu$  is unknown, where

$$c_n = \sqrt{\frac{n}{2}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n+1}{2})}.$$

**Example 5.2.1** (Continuation). Find the Information Lower Bounds for the variances of unbiased estimators of  $\sigma^2$  and  $\sigma$  when  $\mu$  is known and when  $\mu$  is unknown. Do the variances of the UMVUEs obtained above attain these lower bounds?

*Solution.* To find these lower bounds when  $\mu$  is known, we write  $\log f(X; \sigma^2)$  as

$$\ell(X; \theta) = \begin{cases} -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \theta - \frac{1}{2\theta}(X - \mu)^2 & \text{with } \theta = \sigma^2 \\ -\frac{1}{2} \log(2\pi) - \log \theta - \frac{1}{2\theta^2}(X - \mu)^2 & \text{with } \theta = \sigma \end{cases},$$

leading to

$$\dot{\ell}(X; \theta) = \begin{cases} \frac{1}{2\theta^2}[(X - \mu)^2 - \theta] & \text{for } \theta = \sigma^2 \\ \frac{1}{\theta^3}[(X - \mu)^2 - \theta^2] & \text{for } \theta = \sigma \end{cases}.$$

Hence,

$$I(\theta) = E_\theta[\dot{\ell}^2(X; \theta)] = \begin{cases} \frac{1}{2\theta^2} = \frac{1}{2\sigma^4} & \text{for } \theta = \sigma^2 \\ \frac{2}{\theta^2} = \frac{2}{\sigma^2} & \text{for } \theta = \sigma \end{cases},$$

so the information lower bounds for variances of the unbiased estimators of  $\sigma^2$  and  $\sigma$  are, respectively,

$$\frac{1}{nI(\sigma^2)} = \frac{2\sigma^4}{n} \text{ and } \frac{1}{nI(\sigma)} = \frac{\sigma^2}{2n}.$$

When  $\mu$  is unknown, we write  $\log f(X; \mu, \sigma^2)$  as

$$\ell(X; \theta_1, \theta_2) = \begin{cases} -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \theta_2 - \frac{1}{2\theta_2}(X - \theta_1)^2 & \text{with } (\theta_1, \theta_2) = (\mu, \sigma^2) \\ -\frac{1}{2} \log(2\pi) - \log \theta_2 - \frac{1}{2\theta_2^2}(X - \theta_1)^2 & \text{with } (\theta_1, \theta_2) = (\mu, \sigma) \end{cases},$$

leading to

$$\dot{\ell}^\top(X; \theta_1, \theta_2) = \begin{cases} \left( \frac{1}{\theta_2}[X - \theta_1], -\frac{1}{2\theta_2} + \frac{1}{2\theta_2^2}[X - \theta_1]^2 \right) & \text{for } (\theta_1, \theta_2) = (\mu, \sigma^2) \\ \left( \frac{1}{\theta_2^2}[X - \theta_1], -\frac{1}{\theta_2} + \frac{1}{\theta_2^3}[X - \theta_1]^2 \right) & \text{for } (\theta_1, \theta_2) = (\mu, \sigma) \end{cases}.$$

Hence, the information matrix

$$\mathbf{I}(\theta) = E_\theta[\dot{\ell}(X; \theta_1, \theta_2) \dot{\ell}^\top(X; \theta_1, \theta_2)]$$

is obtained as

$$\mathbf{I}(\theta) = \begin{bmatrix} \frac{1}{\theta_2} & 0 \\ 0 & \frac{1}{2\theta_2^2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{bmatrix} \text{ for } (\theta_1, \theta_2) = (\mu, \sigma^2), \text{ and}$$

$$\mathbf{I}(\theta) = \begin{bmatrix} \frac{1}{\theta_2^2} & 0 \\ 0 & \frac{2}{\theta_2^2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{bmatrix} \text{ for } (\theta_1, \theta_2) = (\mu, \sigma),$$

with inverses

$$\mathbf{I}^{-1}(\theta) := \begin{bmatrix} I^{11} & I^{12} \\ I^{21} & I^{22} \end{bmatrix},$$

where

$$I^{22} = \begin{cases} 2\sigma^4 & \text{for } (\theta_1, \theta_2) = (\mu, \sigma^2) \\ \sigma^2/2 & \text{for } (\theta_1, \theta_2) = (\mu, \sigma) \end{cases}.$$

Thus the information lower bounds  $n^{-1}I^{22}$  for variances of unbiased estimators of  $\sigma^2$  and  $\sigma$  are  $2\sigma^4/n$  and  $\sigma^2/(2n)$ , respectively.

We now compare these lower bounds with the variances of the UMVUEs obtained above:

$$\begin{aligned} \text{Var}_\sigma[W] &= n^{-1}\text{Var}_\sigma[(X_i - \mu)^2] \\ &= n^{-1}\left\{\text{E}_\sigma[(X_i - \mu)^4] - \text{E}_\sigma^2[(X_i - \mu)^2]\right\} \\ &= n^{-1}\left[3\sigma^4 - \sigma^4\right] = \frac{2\sigma^4}{n}, \end{aligned}$$

which equals the information lower bound for this case.

$$\begin{aligned} \text{Var}_\sigma[T_0] &= c_n^2 \text{Var}_\sigma[\sqrt{W}] = c_n^2 \left\{\text{E}_\sigma[W] - \text{E}_\sigma^2[\sqrt{W}]\right\} \\ &= c_n^2 \left[\sigma^2 - \left(\frac{\sigma}{c_n}\right)^2\right] = \sigma^2(c_n^2 - 1), \end{aligned}$$

where  $c_n = \sqrt{\frac{n}{2}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n+1}{2})}$ , while the corresponding lower bound is  $\sigma^2/(2n)$ .

Next, we have

$$\text{Var}_\sigma[S^2] = \text{Var}\left[(n-1)^{-1}\sigma^2\chi_{n-1}^2\right] = \frac{2\sigma^4}{n-1},$$

which is greater than the corresponding lower bound  $\frac{2\sigma^4}{n}$ . Finally,

$$\begin{aligned} \text{Var}_\sigma[T_1] &= c_{n-1}^2 \text{Var}_\sigma[S] \\ &= c_{n-1}^2 \left\{\text{E}_\sigma[S^2] - \text{E}_\sigma^2[S]\right\} \\ &= c_{n-1}^2 \left\{\sigma^2 - \left(\frac{\sigma}{c_{n-1}}\right)^2\right\} = \sigma^2(c_{n-1}^2 - 1), \end{aligned}$$

while the corresponding lower bound is  $\sigma^2/(2n)$ .

Thus the variance of the UMVUE of  $\sigma^2$  when  $\mu$  is known (ie,  $\text{Var}_{\sigma^2}[W]$  attains the information lower bound), whereas the other three UMVUEs do not attain their corresponding lower bounds. In particular,

$$\frac{\text{Var}_\sigma[S^2]}{2\sigma^4/n} = \frac{n}{n-1}, \quad \frac{\text{Var}_\sigma[T_0]}{\sigma^2/(2n)} = 2n(c_n^2 - 1), \text{ and}$$

$$\frac{\text{Var}_\sigma[T_1]}{\sigma^2/(2n)} = 2n(c_{n-1}^2 - 1).$$

However, all three of these ratios are nearly equal to 1 for moderately large  $n$ . For example, for  $n = 20$ ,

$$\frac{n}{n-1} = 1.0526, \quad 2n(c_n^2 - 1) = 1.0144, \text{ and } 2n(c_{n-1}^2 - 1) = 1.0630.$$

### 5.3 Equivariance

Like unbiasedness, equivariance is another property which narrows down the class of estimators to only those with constant risk functions, allowing us to look for the best among them.

We first briefly discuss equivariance in the general context of a decision problem specified by

$$(\mathcal{X}, \mathcal{B}, \{P_\theta, \theta \in \Theta\}), A, L$$

in the light of a group  $\mathcal{G}$  of one-to-one transformations of  $\mathcal{X} \xrightarrow{\text{onto}} \mathcal{X}$ , the group operation being composition, that is,

$$g_2 \cdot g_1(x) = g_2[g_1(x)].$$

The identity transformation is  $e(x) = x$  for all  $x$ , the inverse of  $g$  is the usual inverse transformation  $g^{-1}$  defined by  $g^{-1}[g(x)] = x$  for all  $x$ . We assume that the family  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  is identifiable (ie, if  $\theta \neq \theta'$ , then  $P_\theta \neq P_{\theta'}$ ).

The family of probabilities  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  is said to be *invariant under  $\mathcal{G}$*  if for every  $g \in \mathcal{G}$  and for every  $\theta \in \Theta$ , there exists  $\theta' \in \Theta$  such that if  $X \sim P_\theta$ , then  $g(X) \sim P_{\theta'}$ . Since the family  $\{P_\theta, \theta \in \Theta\}$  is identifiable, such  $\theta'$  is unique and we denote it by  $\bar{g}(\theta)$ . Thus for any  $B$ ,

$$P_\theta[g^{-1}(B)] = P_\theta[X \in g^{-1}(B)] = P_\theta[g(X) \in B] = P_{\bar{g}(\theta)}[B].$$

Let  $\bar{\mathcal{G}} = \{\bar{g}: g \in \mathcal{G}\}$ . Then  $\bar{\mathcal{G}}$  is a group of transformations from  $\Theta \rightarrow \Theta$ , all  $\bar{g} \in \bar{\mathcal{G}}$  being one-to-one and onto. [To check that  $\bar{\mathcal{G}}$  is a group, verify (i)  $\bar{g}_2 \cdot \bar{g}_1 = \overline{g_2 \cdot g_1}$ , (ii)  $\bar{e}$  is the identity element of  $\bar{\mathcal{G}}$ , and (iii)  $\bar{g}^{-1} = \overline{g^{-1}}$ .]

**Definition 5.3.1.** A statistical decision problem specified by

$$(\mathcal{X}, \mathcal{B}, \{P_\theta, \theta \in \Theta\}), A, L$$

is said to be invariant under a group  $\mathcal{G}$  of transformations from  $\mathcal{X} \rightarrow \mathcal{X}$  if

- (i)  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  is invariant under  $\mathcal{G}$  with the associated group of transformations  $\bar{\mathcal{G}}$  from  $\Theta \rightarrow \Theta$ , and

- (ii) for every  $\bar{g} \in \tilde{\mathcal{G}}$  and for every  $a \in A$ , there is a unique  $a' \in A$  denoted by  $a' = \tilde{g}(a)$  such that for all  $\theta \in \Theta$

$$L(\theta, a) = L(\bar{g}(\theta), a') = L(\bar{g}(\theta), \tilde{g}(a)).$$

*Remark 5.3.1.* If for a given  $\bar{g} \in \tilde{\mathcal{G}}$  and  $a \in A$ , (ii) is satisfied by more than one  $a'$ , then we can remove all such candidates except one and call this  $a' = \tilde{g}(a)$ , thereby achieving uniqueness of  $a'$  as specified in (ii).

Let  $\tilde{\mathcal{G}} = \{\tilde{g}: g \in \mathcal{G}\}$ . Then  $\tilde{\mathcal{G}}$  is a group of transformations from  $A \rightarrow A$ . All  $\tilde{g}$  are one-to-one and onto. [Verify that (i)  $\tilde{g}_2 \cdot \tilde{g}_1 = \widetilde{g_2 \cdot g_1}$ , (ii)  $\tilde{e}$  is the identity transformation of  $\tilde{\mathcal{G}}$ , and (iii)  $\tilde{g}^{-1} = \widetilde{g^{-1}}$ .]

**Definition 5.3.2.** If a decision problem is invariant under  $\mathcal{G}$ , then

- (i) a nonrandomized decision rule  $d$  is equivariant under  $\mathcal{G}$  if for all  $x \in \mathcal{X}$  and  $g \in \mathcal{G}$ ,  $d(g(x)) = \tilde{g}(d(x))$ , and
- (ii) a behavioral decision rule  $\delta$  is equivariant under  $\mathcal{G}$  if  $\delta(S|g(x)) = \delta(\tilde{g}^{-1}(S)|x)$  for all “events”  $S \subset A$ , for all  $x \in \mathcal{X}$  and  $g \in \mathcal{G}$ .

**Definition 5.3.3.**

- (a) For  $\theta_1, \theta_2 \in \Theta$ , we say  $\theta_1 \equiv \theta_2$  if there exists  $\bar{g} \in \tilde{\mathcal{G}}$  such that  $\theta_2 = \bar{g}(\theta_1)$ ,  $\equiv$  being an equivalence relation;
- (b) let  $\Theta^*$  denote the collection of equivalence classes of  $\Theta$ ; and
- (c) the equivalence classes are called orbits,  $\Theta^*$  being the set of all orbits.

For the rest of this section, we consider only nonrandomized decision rules.

**Theorem 5.3.1.** *For any (nonrandomized) equivariant decision rule  $d$ , the risk function is constant on each orbit.*

*Proof.* Note that

$$\begin{aligned} R(\theta, d) &= E_\theta[L(\theta, d(X))] = E_\theta[L(\bar{g}(\theta), \tilde{g}(d(X)))] \\ &= E_\theta[L(\bar{g}(\theta), d(g(X)))] = E_{\bar{g}(\theta)}[L(\bar{g}(\theta), d(Y))] = R(\bar{g}(\theta), d), \end{aligned}$$

using (i)  $L(\theta, a) = L(\bar{g}(\theta), \tilde{g}(a))$ , (ii)  $\tilde{g}(d(x)) = d(g(x))$  since  $d$  is equivariant, and (iii)  $X \sim P_\theta$  implies  $Y = g(X) \sim P_{\bar{g}(\theta)}$ .  $\square$

This theorem shows how equivariance simplifies the search for an optimal rule. If there is a single orbit, then there is an optimal equivariant rule, which is called the Minimum Risk Equivariant (MRE) estimator.

### 5.3.1 Location Equivariance in a Location Family

Suppose  $\mathbf{X} = (X_1, \dots, X_n)$  has pdf  $f(\mathbf{x}; \theta) = f(x_1 - \theta, \dots, x_n - \theta)$  where  $\theta$  is the unknown location parameter to be estimated.

Let  $\mathbf{Y} = g_c(\mathbf{X}) = (X_1 + c, \dots, X_n + c) := (\mathbf{X} + c\mathbf{1})$  where  $\mathbf{1}^\top = (1, \dots, 1)$ , and let  $\bar{g}_c(\theta) = \theta + c$ . Then  $\mathbf{Y}$  has pdf

$$f([y_1 - c] - \theta, \dots, [y_n - c] - \theta) = f(\mathbf{y} - [\theta + c]\mathbf{1}) = f(\mathbf{y}; \bar{g}_c(\theta)).$$

For  $a \in \mathbb{R}$ , taking  $\tilde{g}_c(a) = a + c$ , suppose

$$L(\bar{g}_c(\theta), \tilde{g}_c(a)) = L(\theta + c, a + c) = L(\theta, a) \quad \text{for all } \theta, a, c.$$

Then with  $c = -\theta$ ,

$$L(\theta, a) = L(0, a - \theta) := \rho(a - \theta).$$

Thus the problem is invariant under  $\mathcal{G} = \{g_c, c \in \mathbb{R}\}$  and a location equivariant rule  $d$  must satisfy:  $d(\mathbf{x} + c\mathbf{1}) = d(\mathbf{x}) + c$ . All the usual estimators of a location parameter such as mean, median, etc., are location equivariant.

In this set-up, there is only one orbit in  $\Theta = \mathbb{R}$  and therefore, the risk function is constant in  $\theta$  for equivariant rules  $d$  satisfying  $d(\mathbf{x} + c\mathbf{1}) = d(\mathbf{x}) + c$ . Indeed, for such a rule,

$$\begin{aligned} R(\theta, d) &= \int L(\theta, d(\mathbf{x}))f(\mathbf{x} - \theta\mathbf{1}) d\mathbf{x} \\ &= \int L(\theta, d(\mathbf{y} + \theta\mathbf{1}))f(\mathbf{y}) d\mathbf{y} \\ &= \int L(\theta, d(\mathbf{y}) + \theta)f(\mathbf{y}) d\mathbf{y} \\ &= \int L(0, d(\mathbf{y}))f(\mathbf{y}) d\mathbf{y} = R(0, d). \end{aligned}$$

We now find the MRE estimator of  $\theta$ . We first characterize a location equivariant estimator.

**Lemma 5.3.1.**  *$d$  is location equivariant iff  $d(\mathbf{x}) = d_0(\mathbf{x}) + u(\mathbf{x})$  where  $d_0$  is an arbitrary location equivariant estimator and  $u(\mathbf{x} + c\mathbf{1}) = u(\mathbf{x})$  for all  $\mathbf{x}$  and  $c \in \mathbb{R}$ .*

*Proof.* If  $d$  is of the above form, then

$$\begin{aligned} d(\mathbf{x} + c\mathbf{1}) &= d_0(\mathbf{x} + c\mathbf{1}) + u(\mathbf{x} + c\mathbf{1}) \\ &= \{d_0(\mathbf{x}) + c\} + u(\mathbf{x}) \\ &= \{d_0(\mathbf{x}) + u(\mathbf{x})\} + c = d(\mathbf{x}) + c \end{aligned}$$

showing that  $d$  is equivariant. Conversely, if  $d$  is equivariant, then taking  $u(\mathbf{x}) = d(\mathbf{x}) - d_0(\mathbf{x})$ , we have  $d(\mathbf{x}) = d_0(\mathbf{x}) + u(\mathbf{x})$  and

$$\begin{aligned} u(\mathbf{x} + c\mathbf{1}) &= d(\mathbf{x} + c\mathbf{1}) - d_0(\mathbf{x} + c\mathbf{1}) \\ &= \{d(\mathbf{x}) + c\} - \{d_0(\mathbf{x}) + c\} \\ &= d(\mathbf{x}) - d_0(\mathbf{x}) = u(\mathbf{x}). \end{aligned}$$

□

To find the MRE estimator,

1. For  $n = 1$ , take  $d_0(x) = x$  which is location equivariant, and

$$u(x + c) = u(x) \quad \text{for all } x \text{ and } c \in \mathbb{R}$$

implies  $u(x) = \text{constant} := -b$ . Hence, all location equivariant estimators are of the form:  $d(x) = x - b := d_b(x)$ . Now the MRE estimator is obtained by minimizing

$$R(0, d_b) = E_0[\rho(X - b)]$$

with respect to  $b$ .

2. For  $n \geq 2$ , let  $\mathbf{y} = (y_1, \dots, y_{n-1}) = (x_1 - x_n, \dots, x_{n-1} - x_n)$ . Then  $u(\mathbf{x} + c\mathbf{1}) = u(\mathbf{x})$  for all  $\mathbf{x}$  and  $c \in \mathbb{R}$  iff  $u(\mathbf{x}) = v(\mathbf{y})$  for some function  $v$ .

*Proof.*  $u(\mathbf{x}) = v(\mathbf{y})$  means  $u(x_1, \dots, x_n) = v(x_1 - x_n, \dots, x_{n-1} - x_n)$ . Therefore,  $u(\mathbf{x} + c\mathbf{1}) = v(x_1 - x_n, \dots, x_{n-1} - x_n) = u(\mathbf{x})$ . Conversely, if  $u(\mathbf{x} + c\mathbf{1}) = u(\mathbf{x})$ , then with  $c = -x_n$ ,

$$u(\mathbf{x}) = u(x_1 - x_n, \dots, x_{n-1} - x_n, 0) = u(y_1, \dots, y_{n-1}, 0) := v(\mathbf{y}). \quad \square$$

Thus location equivariant estimators are of the form  $d(\mathbf{x}) = d_0(\mathbf{x}) - v(\mathbf{y})$ , so choose  $v$  so that

$$R(0, d) = E_0[\rho(d_0(\mathbf{X}) - v(\mathbf{Y}))] = E_0 E_0[\rho(d_0(\mathbf{X}) - v(\mathbf{Y})) | \mathbf{Y}]$$

is minimum, which is achieved by minimizing  $E_0[\rho(d_0(\mathbf{X}) - v(\mathbf{y})) | \mathbf{Y} = \mathbf{y}]$  for each  $\mathbf{y}$ . Therefore, to find the MRE estimator, choose  $v(\mathbf{y}) = v^*(\mathbf{y})$  for each  $\mathbf{y}$ , so that

$$E_0[\rho(d_0(\mathbf{X}) - v^*(\mathbf{y})) | \mathbf{Y} = \mathbf{y}] \leq E_0[\rho(d_0(\mathbf{X}) - b) | \mathbf{Y} = \mathbf{y}] \quad \text{for all } b.$$

The above development is summarized in the following theorem.

**Theorem 5.3.2.** *For  $n = 1$ , the MRE estimator of a location parameter is given by  $X - b^*$  where*

$$E_0[\rho(X - b^*)] \leq E_0[\rho(X - b)] \quad \text{for all } b.$$

*For  $n \geq 2$ , let  $\mathbf{Y} = (Y_1, \dots, Y_{n-1}) = (X_1 - X_n, \dots, X_{n-1} - X_n)$  and let  $d_0(\mathbf{X})$  be an arbitrary location equivariant estimator with finite risk. Then the MRE estimator is given by  $d_0(\mathbf{X}) - v^*(\mathbf{Y})$  where for each  $\mathbf{y}$ ,*

$$E_0[\rho(d_0(\mathbf{X}) - v^*(\mathbf{y})) | \mathbf{Y} = \mathbf{y}] \leq E_0[\rho(d_0(\mathbf{X}) - b) | \mathbf{Y} = \mathbf{y}] \quad \text{for all } b.$$

In particular, for squared-error loss  $L(\theta, a) = (a - \theta)^2$ , the MRE estimator is:

$$\begin{aligned} d^*(x) &= x - E_0[X] \text{ for } n = 1 \text{ and} \\ d^*(\mathbf{x}) &= d_0(\mathbf{x}) - E_0[d_0(\mathbf{X}) | \mathbf{Y} = \mathbf{y}] \text{ for } n \geq 2, \end{aligned}$$

because  $E_0[(X - b)^2]$  and  $E_0[\{d_0(\mathbf{X}) - b\}^2 | \mathbf{Y} = \mathbf{y}]$  are minimized at  $b = E_0[X]$  and  $b = E_0[d_0(\mathbf{X}) | \mathbf{Y} = \mathbf{y}]$ , respectively.

Likewise, for absolute error loss  $L(\theta, a) = |a - \theta|$ , the MRE estimator is:

$$\begin{aligned} d^*(x) &= x - \text{median}_0[X] \text{ for } n = 1 \text{ and} \\ d^*(\mathbf{x}) &= d_0(\mathbf{x}) - \text{median}_0[d_0(\mathbf{X}) | \mathbf{Y} = \mathbf{y}] \text{ for } n \geq 2. \end{aligned}$$

Suppose  $L(\theta, a) = \rho(a - \theta)$  where  $\rho$  is convex and even. Suppose further that  $f$  is symmetric about 0. Then for  $n = 1$ , the MRE estimator is  $d^*(x) = x$ . This is because

$$\begin{aligned} E_0[\rho(X - b)] &= \int_{-\infty}^0 \rho(x - b)f(x) dx + \int_0^\infty \rho(x - b)f(x) dx \\ &= \int_0^\infty \rho(-x - b)f(-x) dx + \int_0^\infty \rho(x - b)f(x) dx \\ &= \int_0^\infty [\rho(x + b) + \rho(x - b)]f(x) dx, \end{aligned}$$

so that

$$E_0[\rho(X - b)] - E_0[\rho(X)] = 2 \int_0^\infty \left[ \frac{1}{2} \{\rho(x + b) + \rho(x - b)\} - \rho(x) \right] f(x) dx \geq 0$$

for all  $b$ .

**Example 5.3.1.** Let  $(X_1, \dots, X_n)$  be a random sample from  $\mathcal{N}(\theta, \sigma^2)$  where  $\sigma^2$  is known. Take  $d_0(\mathbf{X}) = \bar{X}$  which is complete sufficient, while

$$\mathbf{Y} = (X_1 - X_n, \dots, X_{n-1} - X_n)$$

is an *ancillary statistic* (ie, its distribution does not depend on  $\theta$ ). By Basu's [Theorem 5.1.6](#), it follows that  $d_0(\mathbf{X}) = \bar{X}$  is independent of  $\mathbf{Y}$ . Hence, for each  $\mathbf{y}$ ,

$$E_0[\rho(d_0(\mathbf{X}) - b) | \mathbf{Y} = \mathbf{y}] = E_0[\rho(\bar{X} - b) | \mathbf{Y} = \mathbf{y}] = E_0[\rho(\bar{X} - b)].$$

For  $\theta = 0$ , the pdf of  $\bar{X}$  is symmetric about 0, so that for all convex and even  $\rho$ ,  $E_0[\rho(\bar{X} - b)]$  is minimized for  $b = 0$ . Thus the MRE estimator is  $d^*(\mathbf{X}) = d_0(\mathbf{X}) - 0 = \bar{X}$ .

**Theorem 5.3.3.** Under squared-error loss, the MRE estimator

$$d^*(\mathbf{x}) = d_0(\mathbf{x}) - E_0[d_0(\mathbf{X}) | \mathbf{Y} = \mathbf{y}]$$

can be expressed as

$$d^*(\mathbf{x}) = \frac{\int_{-\infty}^{\infty} u f(x_1 - u, \dots, x_n - u) du}{\int_{-\infty}^{\infty} f(x_1 - u, \dots, x_n - u) du}.$$

In this form it is known as the Pitman estimator of  $\theta$  [19], which is also the generalized Bayes estimator with respect to the "improper prior having uniform distribution over the entire real line."

*Proof.* Take  $d_0(\mathbf{X}) = X_n$  and compute  $E_0[d_0(\mathbf{X}) | \mathbf{Y}] = E_0[X_n | \mathbf{Y}]$  by augmenting the transformation  $\mathbf{x} \rightarrow \mathbf{y}$  in the following manner:

$$y_1 = x_1 - x_n, \dots, y_{n-1} = x_{n-1} - x_n, \quad y_n = x_n.$$

The Jacobian of this transformation is 1 and the joint pdf of

$$\mathbf{Y}^* = (Y_1, \dots, Y_{n-1}, Y_n)$$

for  $\theta = 0$  is

$$\begin{aligned} f_{\mathbf{Y}^*}(y_1, \dots, y_{n-1}, y_n) &= f_{\bar{X}}(y_1 + y_n, \dots, y_{n-1} + y_n, y_n) \\ &= f(y_1 + y_n, \dots, y_{n-1} + y_n, y_n). \end{aligned}$$

The conditional pdf of  $Y_n = X_n$  given  $\mathbf{Y} = (Y_1, \dots, Y_{n-1})$  is

$$\frac{f_{\mathbf{Y}^*}(y_1, \dots, y_n)}{f_{\mathbf{Y}}(y_1, \dots, y_{n-1})} = \frac{f(y_1 + y_n, \dots, y_{n-1} + y_n, y_n)}{\int_{-\infty}^{\infty} f(y_1 + t, \dots, y_{n-1} + t, t) dt}.$$

Hence, for  $\mathbf{y} = (x_1 - x_n, \dots, x_{n-1} - x_n)$ ,

$$\begin{aligned} E_0[X_n | \mathbf{Y} = \mathbf{y}] &= E_0[Y_n | \mathbf{Y} = \mathbf{y}] \\ &= \frac{\int_{-\infty}^{\infty} tf(y_1 + t, \dots, y_{n-1} + t, t) dt}{\int_{-\infty}^{\infty} f(y_1 + t, \dots, y_{n-1} + t, t) dt} \\ &= \frac{\int_{-\infty}^{\infty} tf(x_1 - x_n + t, \dots, x_{n-1} - x_n + t, t) dt}{\int_{-\infty}^{\infty} f(x_1 - x_n + t, \dots, x_{n-1} - x_n + t, t) dt}, \text{ with } x_n - t = u \\ &= \frac{\int_{-\infty}^{\infty} (x_n - u)f(x_1 - u, \dots, x_{n-1} - u, x_n - u) du}{\int_{-\infty}^{\infty} f(x_1 - u, \dots, x_{n-1} - u, x_n - u) du} \\ &= x_n - \frac{\int_{-\infty}^{\infty} uf(x_1 - u, \dots, x_n - u) du}{\int_{-\infty}^{\infty} f(x_1 - u, \dots, x_n - u) du} \end{aligned}$$

showing that  $d^*(\mathbf{x}) = x_n - E_0[X_n | \mathbf{Y} = \mathbf{y}]$  has the desired form.  $\square$

### 5.3.2 Scale Equivariance in a Scale Family

Suppose  $\mathbf{X} = (X_1, \dots, X_n)$  has pdf  $f(\mathbf{x}; \sigma) = \frac{1}{\sigma^n} f\left(\frac{x_1}{\sigma}, \dots, \frac{x_n}{\sigma}\right)$  where  $\sigma \in \mathbb{R}^+$  is the unknown scale parameter to be estimated.

Let  $\mathbf{Y} = g_c(\mathbf{X}) = c\mathbf{X}$  and  $\tilde{g}_c(\sigma) = c\sigma$ . Then  $\mathbf{Y} = g_c(\mathbf{X})$  has pdf

$$\frac{1}{(c\sigma)^n} f\left(\frac{\mathbf{y}}{c\sigma}\right) = f(\mathbf{y}; \tilde{g}_c(\sigma)),$$

that is,  $\mathbf{Y} = g_c(\mathbf{X})$  has pdf  $f(\mathbf{y}; \tilde{g}_c(\sigma))$ .

For  $a \in A = \mathbb{R}^+$ , taking  $\tilde{g}_c(a) = ca$ , suppose

$$L(\tilde{g}_c(\sigma), \tilde{g}_c(a)) = L(c\sigma, ca) = L(\sigma, a) \quad \text{for all } \sigma, a, c.$$

Then with  $c = 1/\sigma$ ,  $L(\sigma, a) = L(1, a/\sigma) := \rho(a/\sigma)$ .

Thus the problem is invariant under  $\mathcal{G} = \{g_c, c \in \mathbb{R}^+\}$  and an equivariant rule  $d$  must satisfy:  $d(cx) = cd(x)$ . As in the case of location family, here also there is only one orbit in  $\Theta = \mathbb{R}^+$  and therefore, the risk function is constant for equivariant rules.

A rule  $d$  is scale equivariant iff  $d(\mathbf{x}) = u(\mathbf{x})d_0(\mathbf{x})$ , where  $d_0$  is an arbitrary scale equivariant estimator and  $u(cx) = u(x)$  for all  $x$  and  $c \in \mathbb{R}^+$ . The proof is similar to the location case.

We now find the MRE estimator of  $\sigma$ , proceeding as in the location case:

- (i) For  $n = 1$ , take  $d_0(x) = x$  which is scale equivariant. So any scale equivariant  $d$  must be of the form:  $d(x) = u(x)d_0(x) = u(x)x$ , where  $u(cx) = u(x)$  for all  $x$  and  $c \in \mathbb{R}^+$ . Hence,  $u(x) = b^{-1}$  (constant) for all  $x$ . Thus all scale equivariant estimators are of the

form  $d_b(x) = x/b$ . To obtain the MRE estimator, we therefore have to minimize the constant risk of  $d_b$ , namely,

$$R(1, d_b) = R(1, x/b) = E_1[\rho(X/b)]$$

with respect to  $b$ .

- (ii) For  $n \geq 2$ , let  $\mathbf{y} = (y_1, \dots, y_{n-1}) = (x_1/x_n, \dots, x_{n-1}/x_n)$ . Then

$$u(cx) = u(\mathbf{x}) \quad \text{for all } \mathbf{x} \text{ and } c \in \mathbb{R}^+ \Leftrightarrow u(\mathbf{x}) = v(\mathbf{y}) \text{ for some } v,$$

the proof of which is analogous to the corresponding result in the location case. Thus the scale equivariant estimators for  $n \geq 2$  are of the form:  $d(\mathbf{x}) = d_0(\mathbf{x})/v(\mathbf{y})$ . So choose  $v$  so that

$$R(1, d) = E_1[\rho(d_0(\mathbf{X})/v(\mathbf{Y}))] = E_1 E_1[\rho(d_0(\mathbf{X})/v(\mathbf{Y})) | \mathbf{Y}]$$

is minimized. To find the MRE estimator, choose  $v^*(\mathbf{y})$  for each  $\mathbf{y}$  so that

$$E_1[\rho(d_0(\mathbf{X})/v^*(\mathbf{y})) | \mathbf{Y} = \mathbf{y}] \leq E_1[\rho(d_0(\mathbf{X})/b) | \mathbf{Y} = \mathbf{y}] \quad \text{for all } b.$$

The above development is summarized in the following theorem.

**Theorem 5.3.4.** *For  $n = 1$ , the MRE estimator of a scale parameter is given by  $X/b^*$  where*

$$E_1[\rho(X/b^*)] \leq E_1[\rho(X/b)] \quad \text{for all } b.$$

*For  $n \geq 2$ , let  $\mathbf{Y} = (Y_1, \dots, Y_{n-1}) = (X_1/X_n, \dots, X_{n-1}/X_n)$  and let  $d_0(\mathbf{X})$  be an arbitrary scale equivariant estimator with finite risk. Then the MRE estimator is given by  $d_0(\mathbf{X})/v^*(\mathbf{Y})$  where for each  $\mathbf{y}$ ,*

$$E_1[\rho(d_0(\mathbf{X})/v^*(\mathbf{y})) | \mathbf{Y} = \mathbf{y}] \leq E_1[\rho(d_0(\mathbf{X})/b) | \mathbf{Y} = \mathbf{y}] \quad \text{for all } b.$$

## 5.4 Bayesian Estimation Using Conjugate Priors

- (I) In [Section 4.6.3](#) we have introduced the Bayes Principle in the general context of a statistical decision problem described by the triple  $(\{P_\theta, \theta \in \Theta\}, A, L)$ . According to this principle, we choose a decision rule  $d$  for which the risk  $R(\theta, d)$  averaged with respect to a prior distribution  $\tau$  over  $\Theta$  is minimum.
- (II) In [Section 4.6.4](#) the general scheme of finding Bayes rules was outlined, which consists of choosing  $d(x) = d_\tau(x)$  for each observed data  $x$  for which the posterior risk is minimized, that is

$$\int_{\Theta} L(\theta, d_\tau(x)) g(\theta|x) d\theta \leq \int_{\Theta} L(\theta, a) g(\theta|x) d\theta$$

for all  $a \in A$ , where

$$g(\theta|x) = f(x, \theta) \tau(\theta) / \int_{\Theta} f(x, u) \tau(u) du$$

is the posterior pdf of  $\theta$  given  $x$ .

- (III) The method described above is quite straight forward. However, the key step in the method is calculating the posterior distribution  $g(\theta|x)$  given the data, from the prior  $\tau(\theta)$  over  $\Theta$  in a family of pdf's  $\{f(x, \theta), \theta \in \Theta\}$ , which may become very messy. This task becomes very simple in some cases by selecting  $\tau(\theta)$  from a family  $\{\tau(\theta, \alpha), \alpha \in \Omega\}$  which matches with the family  $\{f(x, \theta), \theta \in \Theta\}$  in a manner described below.

Let  $\mathcal{F} = \{f(x, \theta), \theta \in \Theta\}$  be a family of pdf's on  $\mathfrak{X}$  and let  $\tau = \{\tau(\theta, \alpha), \alpha \in \Omega\}$  be a family of priors on  $\Theta$ . Let

$$g_\alpha(\theta|x) = \frac{f(x, \theta)\tau(\theta, \alpha)}{\int_{\Theta} f(x, \theta')\tau(\theta', \alpha) d\theta'}$$

be the posterior pdf of  $\theta$  given  $X = x$  corresponding to the prior  $\tau(\theta, \alpha)$ . If there exists  $\phi: \Omega \times \mathfrak{X} \rightarrow \Omega$  such that  $g_\alpha(\theta|x) = \tau(\theta, \phi(\alpha, x)) = \tau(\theta, \phi_x(\alpha))$ , then  $\tau$  is said to be a family of conjugate priors for  $\theta$ . The calculation of posterior from prior in such a situation can be described by the scheme:

$$\boxed{\text{Prior } \tau(\theta, \alpha)} + \boxed{\text{Data } x} = \boxed{\text{Posterior } \tau(\theta, \phi_x(\alpha))}.$$

Consider the special case of a  $k$ -parameter regular exponential family:

$$\mathcal{F} = \left\{ f(x, \theta) = \exp \left[ \sum_{j=1}^k c_j(\theta) \sum_{i=1}^n T_j(x_i) + \sum_{i=1}^n S(x_i) + nd(\theta) \right] I_A(x), \theta \in \Theta \right\},$$

where  $\Theta = \left\{ \theta \in \mathbb{R}^k : \int_{\mathbb{R}^k} f(x, \theta) dx = 1 \right\}$ . Define

$$\omega(\alpha_1, \dots, \alpha_k, \alpha_{k+1}) = \int_{\theta \in \Theta} \exp \left[ \sum_{j=1}^k \alpha_j c_j(\theta) + \alpha_{k+1} d(\theta) \right] d\theta,$$

and  $\Omega = \{(\alpha_1, \dots, \alpha_k, \alpha_{k+1}) : \omega(\alpha_1, \dots, \alpha_k, \alpha_{k+1}) < \infty\}$ . Let

$$\begin{aligned} \tau(\theta; \alpha_1, \dots, \alpha_{k+1}) &= \exp \left[ \sum_{j=1}^k \alpha_j c_j(\theta) + \alpha_{k+1} d(\theta) - \log \omega(\alpha_1, \dots, \alpha_k, \alpha_{k+1}) \right] \\ &= \frac{\exp \left[ \sum_{j=1}^k \alpha_j c_j(\theta) + \alpha_{k+1} d(\theta) \right]}{\omega(\alpha_1, \dots, \alpha_k, \alpha_{k+1})}, \end{aligned} \tag{3}$$

so that  $\tau(\theta; \alpha_1, \dots, \alpha_k, \alpha_{k+1})$  is a pdf for every  $(\alpha_1, \dots, \alpha_{k+1}) \in \Omega$ . Then

$$\tau = \{ \tau(\theta; \alpha_1, \dots, \alpha_{k+1}), (\alpha_1, \dots, \alpha_k, \alpha_{k+1}) \in \Omega \}$$

is a family of conjugate priors for the exponential family  $\mathcal{F}$ . This is because

$$g_{\alpha_1, \dots, \alpha_{k+1}}(\theta|x) = \tau(\theta, \phi_x(\alpha_1, \dots, \alpha_{k+1})), \tag{3a}$$

where

$$\phi_x(\alpha_1, \dots, \alpha_{k+1}) = \left( \alpha_1 + \sum_{i=1}^n T_1(x_i), \dots, \alpha_k + \sum_{i=1}^n T_k(x_i), \alpha_{k+1} + n \right), \quad (3b)$$

as can be easily verified.

**Example 5.4.1.** Let  $(X_1, \dots, X_n)$  be a random sample from  $Bernoulli(\theta)$ .

Then

$$\begin{aligned} f(\mathbf{x}, \theta) &= \theta^{\sum_i^n x_i} (1-\theta)^{n-\sum_i^n x_i} I_A(\mathbf{x}) \\ &= \exp \left[ \log \left( \frac{\theta}{1-\theta} \right) \sum_{i=1}^n x_i + n \log(1-\theta) \right] I_A(\mathbf{x}), \end{aligned}$$

where  $A = \{0, 1\}^n$  and  $0 < \theta < 1$ . From this we obtain

$$\begin{aligned} \omega(\alpha_1, \alpha_2) &= \int_0^1 \exp \left[ \alpha_1 \log \left( \frac{\theta}{1-\theta} \right) + \alpha_2 \log(1-\theta) \right] d\theta \\ &= \int_0^1 \theta^{\alpha_1} (1-\theta)^{\alpha_2 - \alpha_1} d\theta \\ &= \Gamma(\alpha_1 + 1) \Gamma(\alpha_2 - \alpha_1 + 1) / \Gamma(\alpha_2 + 2), \end{aligned}$$

where  $\alpha_1 + 1 > 0$  and  $\alpha_2 - \alpha_1 + 1 > 0$ . Thus we take

$$\begin{aligned} \tau(\theta; \alpha_1, \alpha_2) &= \frac{\Gamma(\alpha_2 + 2)}{\Gamma(\alpha_1 + 1) \Gamma(\alpha_2 - \alpha_1 + 1)} \exp \left[ \alpha_1 \log \left( \frac{\theta}{1-\theta} \right) + \alpha_2 \log(1-\theta) \right] \\ &= \frac{\Gamma(\alpha_2 + 2)}{\Gamma(\alpha_1 + 1) \Gamma(\alpha_2 - \alpha_1 + 1)} \theta^{\alpha_1} (1-\theta)^{\alpha_2 - 1} \end{aligned}$$

on  $0 < \theta < 1$ , which is  $Beta(\alpha_1 + 1, \alpha_2 - \alpha_1 + 1) := Beta(a, b)$  distribution, where  $a = \alpha_1 + 1$  and  $b = \alpha_2 - \alpha_1 + 1$ . Then

$$\begin{aligned} g_{\alpha_1, \alpha_2}(\theta | \mathbf{x}) &= \tau \left( \theta; \alpha_1 + \sum_{i=1}^n x_i, \alpha_2 + n \right) \\ &= Beta \left( \alpha_1 + \sum_{i=1}^n x_i + 1, \alpha_2 + n - \alpha_1 - \sum_{i=1}^n x_i + 1 \right) \\ &= Beta \left( a + \sum_{i=1}^n x_i, b + n - x_i \right). \end{aligned}$$

Hence, for the squared-error loss, the Bayes estimator with respect to the prior  $\tau(\theta; \alpha_1, \alpha_2)$  is the mean of the distribution  $Beta(a + \sum_{i=1}^n X_i, b + n - \sum_{i=1}^n X_i)$  which is

$$T = T(\mathbf{X}) = \frac{\sum_{i=1}^n X_i + a}{n + a + b}.$$

## 5.5 Methods of Estimation

We have so far discussed the Rao-Blackwell method of finding UMVUEs and two methods of estimation in the general decision theoretic framework following the Bayesian principle and the principle of equivariance. Three other methods, such as,

- (i) the method of maximum likelihood;
- (ii) the method of moments; and
- (iii) the method of minimum  $\chi^2$ ,

will now be discussed, which are applicable in the context of parametric families.

The method of maximum likelihood was introduced by Fisher [20] who also laid the foundation of the theory of estimation and demonstrated the superiority of this method over the method of moments which was widely used for a long time until then. The method of minimum  $\chi^2$  has limited use only in the context of multinomial data to estimate the parameters on which the cell probabilities depend.

### 5.5.1 The Method of Maximum Likelihood

For a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a population with pdf/pmf  $f(x; \theta_1, \dots, \theta_k)$ , the joint pdf/pmf of the data is:

$$f(\mathbf{x}; \theta_1, \dots, \theta_k) = \prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_k),$$

treating  $\mathbf{x}$  as variable and  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$  as fixed. We now treat  $\boldsymbol{\theta}$  as the variable, since we are searching among all possible  $\boldsymbol{\theta} \in \Theta$ , and the data  $\mathbf{x}$  as fixed, because that is what we have in our search for  $\boldsymbol{\theta}$ . This defines the *likelihood function* of  $\boldsymbol{\theta}$  for the data  $\mathbf{x}$  as:

$$L(\boldsymbol{\theta}|\mathbf{x}) = L(\theta_1, \dots, \theta_k|x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta_1, \dots, \theta_k).$$

**Definition 5.5.1.** The maximum likelihood estimate  $\hat{\boldsymbol{\theta}}(\mathbf{x})$  for a given data  $\mathbf{x}$  is defined as the value of  $\boldsymbol{\theta} \in \Theta$  at which  $L(\boldsymbol{\theta}|\mathbf{x})$  attains its maximum, that is

$$\hat{\boldsymbol{\theta}}(\mathbf{x}) = \arg \max_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}|\mathbf{x}) \quad \text{for each } \mathbf{x}.$$

The maximum likelihood estimator (MLE) of  $\boldsymbol{\theta}$  for a random sample  $\mathbf{X}$  is  $\hat{\boldsymbol{\theta}}(\mathbf{X})$ .

The maximization of the likelihood function  $L(\boldsymbol{\theta}|\mathbf{x})$  with respect to  $\boldsymbol{\theta} \in \Theta$  for a given  $\mathbf{x}$  is quite straightforward in many situations. Most often, this maximization is achieved by solving the equation:

$$\frac{\partial}{\partial \theta_i} L(\boldsymbol{\theta}|\mathbf{x}) = 0, \quad i = 1, \dots, k$$

and making sure that the extremum is indeed the maximum. In many situations arising in the context of exponential families, it is easier to work with  $\log L(\theta|x)$ , called the log likelihood.

However, various complications may be encountered in the process of maximization of  $L(\theta|x)$  or  $\log L(\theta|x)$  such as:

- (i)  $\frac{\partial}{\partial \theta_i} \log L(\theta|x) = 0, i = 1, \dots, k$  may have several solutions, in which case we have to search for the global maximum among these solutions,
- (ii) the parameter space  $\Theta$  may be restricted and the solution of the likelihood equations may fall outside the restricted part of  $\Theta$ , in which case we have to modify such a solution appropriately, and
- (iii) the equations  $\frac{\partial}{\partial \theta_i} \log L(\theta|x) = 0, i = 1, \dots, k$  may not have a closed form solution, so the equation may have to be solved iteratively.

Complications like (i) and (ii) will be dealt with in some examples to be discussed in the sequel and (iii) will be discussed in [Chapter 7](#).

The MLEs have an important property known as *invariance*. This means that if  $\eta = g(\theta)$ , then the MLE of  $\eta$  based on  $x$  is  $\hat{\eta}(x) = g(\hat{\theta}(x))$  where  $\hat{\theta}(x)$  is the MLE of  $\theta$  based on  $x$ . The proof of this fact is straightforward if  $g: \Theta \rightarrow \mathcal{H}$  is one-to-one, but it needs a little more care if  $g$  is not one-to-one.

For given data  $x$ , the likelihood function of  $\theta$  is  $L(\theta|x) = f(x, \theta)$ . When  $g$  is one-to-one and onto, then  $\theta = g^{-1}(\eta)$  is well defined for all  $\eta \in \mathcal{H}$  and we define  $L^*(\eta|x) = L(g^{-1}(\eta)|x)$ . In general, when  $g$  is onto but not one-to-one, we define  $L^*(\eta|x) = \sup_{\{\theta: g(\theta)=\eta\}} L(\theta|x)$ .

The proof of  $\hat{\eta}(x) = g(\hat{\theta}(x))$  is now given for the two cases:  $g$  is one-to-one or not.

- (i)  $g$  is one-to-one and onto. If  $\hat{\eta}(x)$  is the MLE of  $\eta$  based on  $x$ , then

$$\begin{aligned} L^*(\hat{\eta}(x)|x) &= \sup_{\eta \in \mathcal{H}} L^*(\eta|x) \iff L(g^{-1}(\hat{\eta}(x))|x) = \sup_{\eta \in \mathcal{H}} L(g^{-1}(\eta)|x), \\ \text{ie, } L(g^{-1}(\hat{\eta}(x))|x) &= \sup_{\theta \in \Theta} L(\theta|x). \end{aligned}$$

But then  $g^{-1}(\hat{\eta}(x)) = \hat{\theta}(x)$ , and so  $\hat{\eta}(x) = g(\hat{\theta}(x))$ .

- (ii)  $g$  is onto but not one-to-one. In this case, if  $\hat{\eta}(x)$  is the MLE of  $\eta$  based on  $x$ , then

$$\sup_{\{\theta: g(\theta)=\hat{\eta}(x)\}} L(\theta|x) = \sup_{\eta \in \mathcal{H}} \sup_{\{\theta: g(\theta)=\eta\}} L(\theta|x) = \sup_{\theta \in \Theta} L(\theta|x),$$

since  $\bigcup_{\eta \in \mathcal{H}} \{\theta: g(\theta) = \eta\} = \Theta$ . Hence  $\hat{\theta}(x) \in \{\theta: g(\theta) = \hat{\eta}(x)\}$ , and so  $g(\hat{\theta}(x)) = \hat{\eta}(x)$ .

The asymptotic properties of the MLEs will be discussed in [Chapter 7](#) along with those of the Likelihood Ratio Tests motivated by the likelihood principle in the context of Hypothesis Testing.

The MLEs are “asymptotically efficient” on one hand, but in finite samples, they may be excessively influenced by a few “outliers.” Methods which modify the MLEs against such weakness will be discussed in [Chapter 10](#).

Here we shall only present a heuristic argument to justify the method of maximum likelihood and in the single-parameter context point out a key feature in the model  $\{f(\mathbf{x}, \theta), \theta \in \mathbb{R}\}$  which determines how well the MLE performs.

To emphasize what happens when the sample size  $n \rightarrow \infty$ , we let the log likelihood  $\log L$  be subscripted by  $n$  and averaged over  $n$  to define

$$L_n^*(\theta) = n^{-1} \log L_n(\theta | X_1, \dots, X_n) = n^{-1} \sum_{i=1}^n \log f(X_i, \theta),$$

which is maximized with respect to  $\theta$ .

Let  $\theta_0$  denote the true value of  $\theta$  and suppose that  $E_{\theta_0}[\log f(X, \theta)]$  is finite for all  $\theta$  in a neighborhood of  $\theta_0$ . Then as  $n \rightarrow \infty$ ,

$$L_n^*(\theta) = n^{-1} \sum_{i=1}^n \log f(X_i, \theta) \xrightarrow{P} E_{\theta_0}[\log f(X, \theta)] := L^*(\theta)$$

at each  $\theta$ . The MLE  $\hat{\theta}_n$  at which  $L_n^*(\theta)$  attains its maximum should, therefore, converge in probability to the value of  $\theta$  at which  $L^*(\theta)$  attains its maximum.

Assume that the family  $\{f(x, \theta), \theta \in \Theta\}$  satisfies the *identifiability condition* which requires that for  $\theta' \neq \theta$ ,  $f(\cdot, \theta')$  is essentially different from  $f(\cdot, \theta)$ . Then for  $\theta \neq \theta_0$ , due to strict concavity of the function  $\log(\cdot)$ , it follows from Jensen's inequality that

$$\begin{aligned} & E_{\theta_0}[\log f(X, \theta) - \log f(X, \theta_0)] \\ &= E_{\theta_0} \left[ \log \frac{f(X, \theta)}{f(X, \theta_0)} \right] < \log E_{\theta_0} \left[ \frac{f(X, \theta)}{f(X, \theta_0)} \right] \\ &= \log \int_{\{x: f(x, \theta_0) > 0\}} [f(x, \theta)/f(x, \theta_0)] f(x, \theta_0) dx \\ &= \log \int_{\{x: f(x, \theta_0) > 0\}} f(x, \theta) dx \leq \log 1 = 0. \end{aligned}$$

Thus for  $\theta \neq \theta_0$ ,  $L^*(\theta) = E_{\theta_0}[\log f(X, \theta)] < E_{\theta_0}[\log f(X, \theta_0)] = L^*(\theta_0)$  (ie, the function  $L^*(\theta)$  has a unique maximum at  $\theta_0$ ).

Now by the WLLN, the graph of  $L_n^*(\theta)$  converges to the graph of  $L^*(\theta)$  at each  $\theta$  with probability tending to 1 as  $n \rightarrow \infty$ . So it seems reasonable that the peaks of the two graphs should approach one another as  $n \rightarrow \infty$  (ie,  $\hat{\theta}_n \rightarrow \theta_0$  in probability).

Now for  $\Theta = \mathbb{R}$ , suppose that the graph of  $L^*(\theta)$  has a large curvature at  $\theta = \theta_0$ . Then it falls off sharply as  $\theta$  moves away from its peak at  $\theta_0$ . This increases the tendency of the peak of  $L_n^*(\theta)$  to stay close to the peak of  $L^*(\theta)$ . The geometry of the graph of  $L^*(\theta)$  and the convergence property of  $L_n^*(\theta)$  together suggest that

1.  $\hat{\theta}_n \xrightarrow{P} \theta_0$  as  $n \rightarrow \infty$ , and
2. the precision of  $\hat{\theta}_n$  in estimating  $\theta_0$  increases with the curvature of  $L^*(\theta)$  at  $\theta = \theta_0$ .

The curvature of  $L^*(\theta)$  at  $\theta = \theta_0$  is the Fisher-information  $I(\theta_0)$  of the family  $\{f(x, \theta), \theta \in \mathbb{R}\}$  at  $\theta_0$  (defined in Eq. (2)):

$$\begin{aligned} -\frac{d^2}{d\theta^2}L^*(\theta)\Big|_{\theta=\theta_0} &= -\frac{d^2}{d\theta^2}\int\{\log f(x, \theta)\}f(x, \theta_0) dx\Big|_{\theta=\theta_0} \\ &= -\int\left(\frac{\partial^2 \log f(x, \theta_0)}{\partial \theta^2}\right)f(x, \theta_0) dx \\ &= -E_{\theta_0}\left[\frac{\partial^2 \log f(X, \theta_0)}{\partial \theta^2}\right] \\ &= E_{\theta_0}\left[\left(\frac{\partial \log f(X, \theta_0)}{\partial \theta}\right)^2\right] = I(\theta_0) \end{aligned}$$

under regularity conditions. Note that the Cramér-Rao lower bound is also determined by this curvature, which is an intrinsic measure of how well  $\theta_0$  can be estimated.

We now pursue this heuristic argument to examine the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_n - \theta_0)$ . Assume that the MLE  $\hat{\theta}_n$  satisfies

$$n^{-1}\sum_{i=1}^n \frac{\partial \log f(X_i, \hat{\theta}_n)}{\partial \theta} = 0, \quad (4)$$

and expand the left-hand side of this equation around  $\theta_0$  to see that

$$0 = n^{-1}\sum_{i=1}^n \frac{\partial \log f(X_i, \theta_0)}{\partial \theta} + (\hat{\theta}_n - \theta_0)n^{-1}\sum_{i=1}^n \frac{\partial^2 \log f(X_i, \hat{\theta}_n)}{\partial \theta^2}$$

with  $\hat{\theta}_n$  lying between  $\hat{\theta}_n$  and  $\theta_0$ . Thus

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{n^{-1/2}\sum_{i=1}^n \frac{\partial \log f(X_i, \theta_0)}{\partial \theta}}{n^{-1}\sum_{i=1}^n \frac{\partial^2 \log f(X_i, \hat{\theta}_n)}{\partial \theta^2}},$$

of which the numerator  $\xrightarrow{\mathcal{L}} N(0, I(\theta_0))$  by the CLT and the denominator  $\xrightarrow{P} I(\theta_0)$  by the WLLN. Hence,  $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta_0))$  by Slutsky's Theorem.

*Remark 5.5.1.* The condition that the MLE  $\hat{\theta}_n$  satisfies Eq. (4) is essential for the asymptotic normality of  $\sqrt{n}(\hat{\theta}_n - \theta_0)$ . To see what happens otherwise, let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $Uniform(0, \theta)$ . Then the MLE  $\hat{\theta}_n$  of  $\theta$  is  $X_{n:n} = \max(X_1, \dots, X_n)$  and  $P_{\theta_0}[\hat{\theta}_n \leq t] = (t/\theta_0)^n$  for  $0 < t < \theta_0$  and  $= 1$  for  $t \geq \theta_0$ .

In Chapter 7, the results indicated above by heuristic arguments will be proved more systematically.

### 5.5.2 The Method of Moments

In a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from a population with pdf/pmf  $f(x; \theta_1, \dots, \theta_k)$ , let

$$\mu_j(\theta_1, \dots, \theta_k) = E_{\theta_1, \dots, \theta_k}[X_i^j] \text{ and } m_{nj} = n^{-1} \sum_{i=1}^n X_i^j, \quad j = 1, \dots, k$$

denote, respectively, the  $j$ th population and sample moments. The method of moments estimators (MOME) of  $\theta_1, \dots, \theta_k$  are obtained by equating each of the first  $k$  population moments to the corresponding sample moments and solving the system of equations:

$$\mu_j(\theta_1, \dots, \theta_k) = m_{nj}, \quad j = 1, \dots, k$$

for  $\theta_1, \dots, \theta_k$ . Suppose these equations have unique solutions and let

$$\hat{\theta}_{nj} = g_j(m_{n1}, \dots, m_{nk}), \quad j = 1, \dots, k$$

denote the solutions of these equations. If the population distribution has finite  $2k$  moments, then  $\sqrt{n}(m_{nj} - \mu_j)$ ,  $j = 1, \dots, k$  will be asymptotically jointly normal and consequently,  $\sqrt{n}(\hat{\theta}_{nj} - \theta_j)$ ,  $j = 1, \dots, k$  will also be asymptotically jointly normal, provided that the functions  $g_1, \dots, g_k$  are sufficiently smooth.

**Example 5.5.1.**  $X_1, \dots, X_n$  are iid  $Bernoulli(\theta)$ .

Here  $k = 1$ ,  $\mu(\theta) = \theta$  and  $m_n = n^{-1} \sum_1^n X_i = \bar{X}_n$ . So we equate  $\mu(\theta) = m_n$  to obtain  $\hat{\theta}_n = \bar{X}_n$  which is the same as the MLE.

**Example 5.5.2.**  $X_1, \dots, X_n$  are iid  $Uniform(\theta)$ .

Here also  $k = 1$ ,  $\mu(\theta) = \theta/2$  and  $m_n = \bar{X}_n$ ; so we equate  $\mu(\theta) = m_n$  to obtain  $\hat{\theta}_n = 2\bar{X}_n$ . Since  $E_\theta[X_i] = \theta/2$  and  $Var_\theta[X_i] = \theta^2/12$ , we have  $E_\theta[2\bar{X}_n] = \theta$  and  $Var_\theta[2\bar{X}_n] = \theta^2/3$ .

Hence,  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} \mathcal{N}(0, \theta^2/3)$  which is quite different from the behavior of the MLE. It is left as an exercise to compare the rates at which  $|\hat{\theta}_n - \theta| < \epsilon$  converge to 1 for the MOME and the MLE and to note that the convergence is much faster for the MLE.

*Remark 5.5.2.* Two major drawbacks of the method of moments are:

1. The MOMEs are less efficient than the MLEs in many situations as illustrated in [Example 5.5.2](#).
2. The method is inapplicable if the required number of moments do not exist, as in the case of estimating the median of  $Cauchy(\theta)$ .

### 5.5.3 The Method of Minimum $\chi^2$

Consider a multinomial distribution in  $m$  classes with probability  $\pi_j(\theta)$ ,  $1 \leq j \leq m$  for the  $j$ th class where  $\pi_1(\cdot), \dots, \pi_m(\cdot)$  are known functions of an unknown  $k$ -dimensional parameter vector  $\theta$ , satisfying  $\pi_j(\theta) > 0$  for all  $j$  and  $\sum_{j=1}^m \pi_j(\theta) = 1$ . Let  $n_1, \dots, n_m$  denote the observed frequencies in the  $m$  classes in a random sample of size  $n$ .

In order to put this in a framework suitable for pursuing maximum likelihood estimation of  $\theta$ , let  $e_j$  denote the  $m$ -dimensional vector with 1 for the  $j$ th coordinate and 0

for all other coordinates. Then  $X_i = e_j$  means that in the multinomial sampling, the  $i$ th observation is in the  $j$ th class. The random vectors  $X_1, \dots, X_n$  are iid with probabilities  $P_{\theta}[X_i = e_j] = f(e_j, \theta) = \pi_j(\theta)$ , that is, the pmf of  $X_i$  is

$$f(x, \theta) = \begin{cases} \pi_j(\theta) & \text{if } x = e_j, \quad j = 1, \dots, m \\ 0 & \text{otherwise} \end{cases}$$

We now have for  $x = e_j$ ,

$$\frac{\partial \log f(x, \theta)}{\partial \theta_r} = \frac{\partial \log \pi_j(\theta)}{\partial \theta_r} = \frac{1}{\pi_j(\theta)} \frac{\partial \pi_j(\theta)}{\partial \theta_r},$$

and

$$\begin{aligned} S_{nr}(\theta) &= \sum_{i=1}^n \frac{\partial \log f(X_i, \theta)}{\partial \theta_r} = \sum_{i=1}^n \sum_{j=1}^m I[X_i = e_j] \frac{\partial \log \pi_j(\theta)}{\partial \theta_r} \\ &= \sum_{j=1}^m \left\{ \sum_{i=1}^n I[X_i = e_j] \right\} \frac{1}{\pi_j(\theta)} \frac{\partial \pi_j(\theta)}{\partial \theta_r} = \sum_{j=1}^m \frac{n_j}{\pi_j(\theta)} \frac{\partial \pi_j(\theta)}{\partial \theta_r}, \end{aligned}$$

since  $\sum_{i=1}^n I[X_i = e_j] = \{\text{number of observations } X_i = e_j\} = n_j$ .

The MLEs of  $\theta_1, \dots, \theta_k$  are obtained by solving the likelihood equations:  $S_{nr}(\hat{\theta}_n) = 0$ ,  $1 \leq r \leq k$ , that is

$$\sum_{j=1}^m \frac{n_j}{\pi_j(\hat{\theta}_n)} \frac{\partial \pi_j(\hat{\theta}_n)}{\partial \theta_r} = 0, \quad 1 \leq r \leq k.$$

These equations do not have a closed form solution in a typical problem, so they have to be solved iteratively. This issue will be discussed in [Chapter 7](#).

For frequency data described above, one can try to estimate  $\theta$  by minimizing a measure of discrepancy between  $(n\pi_1(\theta_n^*), \dots, n\pi_m(\theta_n^*))$  which are the *expected frequencies* for  $(\pi_1(\theta_n^*), \dots, \pi_m(\theta_n^*))$  and the *observed frequencies*  $(n_1, \dots, n_m)$  given by

$$\chi^2(\theta_n^*) = \sum_{j=1}^m \frac{(n_j - n\pi_j(\theta_n^*))^2}{n\pi_j(\theta_n^*)} = \sum_{j=1}^m \frac{n_j^2}{n\pi_j(\theta_n^*)} - n$$

with respect to  $\theta_n^*$ . Such a procedure is known as the *method of minimum  $\chi^2$* . The minimizing equations for this method are

$$\sum_{j=1}^m \frac{n_j^2}{\pi_j^2(\theta_n^*)} \frac{\partial \pi_j(\theta_n^*)}{\partial \theta_r} = 0, \quad 1 \leq r \leq k,$$

whereas the likelihood equations are

$$\sum_{j=1}^m \frac{n_j}{\pi_j(\hat{\theta}_n)} \frac{\partial \pi_j(\hat{\theta}_n)}{\partial \theta_r} = 0, \quad 1 \leq r \leq k.$$

The two sets of equations can be shown to be asymptotically equivalent so  $\hat{\theta}_n^*$  has the same asymptotic distribution as that of  $\hat{\theta}_n$ . Thus the estimator  $\hat{\theta}_n^*$  obtained by the method of minimum  $\chi^2$  is asymptotically efficient.

However, for moderate sample sizes, the MLE  $\hat{\theta}_n$  and the minimum  $\chi^2$  estimator  $\hat{\theta}_n^*$  have been known to perform quite differently. The asymptotics being discussed here do not take into account what happens in respect of terms which are of order of magnitude less than  $1/\sqrt{n}$  in probability. What happens beyond this order of magnitude can cause differences in the precision of two estimates which are both “asymptotically efficient” as measured by “first-order efficiency.” These concerns have led to several approaches to develop a measure of “second-order efficiency” [21].

## Exercises

For all the problems below, when there are observations  $X_1, \dots, X_n$ , it is understood that

$$\bar{X}_n = n^{-1} \sum_{i=1}^n X_i, \quad s_n^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

$$X_{n:1} = \min(X_1, \dots, X_n), \quad X_{n:n} = \max(X_1, \dots, X_n).$$

- 5.1.** Let  $X_1, \dots, X_n$  be a random sample from a population with pdf

$$f(x, \theta) = \theta \exp(-\theta x), \quad x > 0, \quad \theta > 0.$$

Find the UMVUE of  $\gamma(\theta) = \exp(-\theta) = P_\theta[X_i > 1]$ .

- 5.2.** Let  $X_1, \dots, X_n$  be a random sample from

$$f(x, \theta) = \exp(-(x - \theta)), \quad x > \theta.$$

- (a) Show that  $T = \min(X_1, \dots, X_n)$  is a complete sufficient statistic for  $\theta$ . [For completeness, you may restrict attention to continuous functions  $\phi$  with  $E_\theta[\phi(T)] = 0$  for all  $\theta$ .]

- (b) Calculate  $E_\theta[T]$  and find a UMVUE of  $\theta$ .

- 5.3.** Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ . Find the UMVUE of  $\sigma$ .

- 5.4.** Let  $X_1, \dots, X_n$  be a random sample from  $Unif(\theta_1, \theta_2)$ . We want to estimate the mean  $\gamma = (\theta_1 + \theta_2)/2$ .

- (a) Show that  $(X_{n:1}, X_{n:n})$  is complete and sufficient for  $(\theta_1, \theta_2)$ , and  $\bar{X}_n$  is an unbiased estimate of  $\gamma$ .

- (b) Find the UMVUE of  $\gamma$ .

- 5.5.** Let  $X_1, \dots, X_n$  be a random sample from  $N(\theta, 1)$ . Find the UMVUE of  $P_\theta[X \geq 0] = \Phi(\theta)$ , where  $\Phi$  is the cdf of the standard normal distribution. [Hint:  $I_{[0, \infty)}(X_1)$  is unbiased for  $\Phi(\theta)$  and  $(X_1, \bar{X}_n)$  is bivariate normal.]

- 5.6.** Let  $X \sim Bin(n, \theta)$ . Show that  $X(n - X)/\{n(n - 1)\}$  is UMVUE of  $\gamma = \theta(1 - \theta)$ .

- 5.7.** Suppose that  $T_1$  and  $T_2$  are two UMVUEs of  $g(\theta)$  with finite variance. Show that  $T_1 = T_2$ . [Hint:  $(T_1 + T_2)/2$  is also unbiased; use the correlation inequality.]

- 5.8.** Show that for an exponential family (in natural form)

$$g(x, \eta) = \exp(\eta T(x) + d_0(\eta) + S(x)) I_A(x),$$

the Fisher-information is  $I_g(\eta) = -d_0''(\eta)$ .

- 5.9.** Consider the exponential family

$$f(x, \theta) = \exp(C(\theta)T(x) + d(\theta) + S(x)) I_A(x),$$

where  $C$  is one-to-one and twice differentiable. Then we can reparametrize  $f(x, \theta)$  by letting  $\eta = C(\theta)$  and rewrite  $f(x, \theta)$  in the form  $g(x, \eta)$  given in Exercise 5.8 above. Show that the Fisher-information in the family  $\{f(x, \theta)\}$  is

$$I_f(\theta) = -d_0''(C(\theta)) = -d_0''(\eta)|_{\eta=C(\theta)} \{C'(\theta)\}^2 = I_g(C(\theta)) \{C'(\theta)\}^2.$$

- 5.10.** Suppose that  $X_1, \dots, X_n$  is a random sample from the pdf

$$f(x, \theta) = cx^{c-1}\theta \exp(-\theta x^c), \quad x > 0, \quad c > 0 \text{ (known)}, \quad \theta > 0 \text{ (unknown)}.$$

Show that  $T = n^{-1} \sum_{i=1}^n X_i^c$  is the UMVUE of  $1/\theta$ .

- 5.11.** Let  $X$  be a random sample from the Cauchy distribution with pdf

$$f(x, \theta) = \frac{1}{\pi \{1 + (x - \theta)^2\}}, \quad -\infty < x < \infty.$$

We want to estimate the median of  $\theta$  under the loss function  $L(\theta, a) = \rho(a - \theta)$ , where  $\rho(t) = I_{(c, \infty)}(|t|)$ . Find the MRE estimator of  $\theta$  based on  $X$ .

- 5.12.** Let  $X_1, \dots, X_n$  be a random sample from  $\text{Unif}(\theta - 1/2, \theta + 1/2)$ . Find the MRE estimator of  $\theta$  under the loss function  $L(\theta, a) = |a - \theta|$ .

- 5.13.** Let  $X_1, \dots, X_n$  be a random sample from the half-normal distribution with pdf

$$f(x, \theta) = \sqrt{2/\pi} \exp[-(x - \theta)^2/2] I_{(0, \infty)}(x).$$

Show that the Pitman estimator of  $\theta$  is

$$d(\mathbf{x}) = \bar{X}_n - \frac{\exp[-n(X_{n:1} - \bar{X}_n)^2/2]}{\sqrt{2n\pi} \Phi(\sqrt{n}(X_{n:1} - \bar{X}_n)^2)},$$

where  $\Phi$  is the cdf of  $N(0, 1)$ .

- 5.14.** Let  $(X_1, \dots, X_n)$  have a joint distribution with pdf

$$f(x_1, \dots, x_n; \theta) = \theta^{-n} g(x_1/\theta, \dots, x_n/\theta)$$

for some function  $g$  which vanishes unless all coordinates are positive, and  $\theta > 0$  is an unknown scale parameter. Suppose that the loss function is  $L(\theta, a) = \rho(a/\theta)$ .

- (a)** Find the MRE estimator of  $\theta$  analogous to the MRE estimator in the location problem. [Take  $\mathbf{Y} = (X_1/X_n, \dots, X_{n-1}/X_n)$ .]

- (b)** For the loss function  $L(\theta, a) = (a/\theta - 1)^2$ , show that the MRE estimator is the following analog of the Pitman estimator:

$$d(\mathbf{x}) = \frac{\int_0^\infty \theta^{-n-2} f(x_1/\theta, \dots, x_n/\theta) d\theta}{\int_0^\infty \theta^{-n-3} f(x_1/\theta, \dots, x_n/\theta) d\theta}.$$

- 5.15.** **(a)** Express the family of beta distributions  $Be(\theta_1, \theta_2)$  as a two-parameter exponential family and find a family of conjugate priors for this family.  
**(b)** Find the prior-to-posterior formula for these conjugate priors. [Use Eqs. (3), (3a), and (3b)]
- 5.16.** Consider the joint pdf  $f(\mathbf{x}, \boldsymbol{\theta})$ , of a random sample  $X_1, \dots, X_n$  from  $N(\theta_1, 1/\theta_2)$ , where  $\boldsymbol{\theta} = (\theta_1, \theta_2)$ .
- (a)** Express  $\mathcal{F}_1 = \{f(\mathbf{x}, \boldsymbol{\theta}): -\infty < \theta_1 < \infty, \theta_2 = 1\}$  as a one-parameter exponential family and find a family of conjugate priors for  $\mathcal{F}_1$ .  
**(b)** Express  $\mathcal{F}_0 = \{f(\mathbf{x}, \boldsymbol{\theta}): \theta_1 = 0, 0 < \theta_2 < \infty\}$  as a one-parameter exponential family and find a family of conjugate priors for  $\mathcal{F}_0$ .  
**(c)** Express  $\mathcal{F}_2 = \{f(\mathbf{x}, \boldsymbol{\theta}): -\infty < \theta_1 < \infty, 0 < \theta_2 < \infty\}$  as a two-parameter exponential family and find a family of conjugate priors for  $\mathcal{F}_2$ .  
**(d)** In each of the cases above, find the prior-to-posterior formula. [Use Eqs. (3), (3a), and (3b).]
- 5.17.** Let  $X_1, \dots, X_n$  be a random sample from a distribution with pdf
- $$f(x, \theta) = \theta x^{-\theta-1} I_{(1, \infty)}(x), \quad \theta > 2.$$
- (a)** Find the method of moment estimator  $\tilde{\theta}_n$  of  $\theta$ .  
**(b)** Determine which of the following is true
- i)  $\tilde{\theta}_n$  is an unbiased estimator of  $\theta$ ,
  - ii)  $\tilde{\theta}_n$  has an upward bias (ie,  $E_\theta(\tilde{\theta}_n) > \theta$ ),
  - iii)  $\tilde{\theta}_n$  has a downward bias. [Write  $\tilde{\theta}_n$  in terms of  $\bar{X}_n - 1$ , find  $E_\theta[\bar{X}_n - 1]$  and use Jensen's inequality.]
- (c)** Find the asymptotic distribution of  $\tilde{\theta}_n$  (suitably normalized) as  $n \rightarrow \infty$ .

- 5.18.** Let  $X_1, \dots, X_n$  be a random sample from a log normal distribution with pdf

$$f(x, \theta) = \frac{1}{\sqrt{2\pi x}} \exp\left[-(\log x - \theta)^2/2\right], \quad x > 0.$$

This means  $Y_i = \log X_i$ ,  $i = 1, \dots, n$ , are iid  $N(\theta, 1)$ , so that  $X_i \stackrel{D}{=} \exp(Z_i + \theta)$  where  $Z_1, \dots, Z_n$  are iid  $N(0, 1)$ . We want to estimate  $\theta$  based on  $(X_1, \dots, X_n)$ .

- (a)** Find the MOME  $\tilde{\theta}_n$  and its asymptotic distribution.  
**(b)** Find the MLE  $\hat{\theta}_n$  and its asymptotic distribution.  
**(c)** Let  $\sigma_1^2$  and  $\sigma_2^2$  be the asymptotic variances of  $\tilde{\theta}_n$  and  $\hat{\theta}_n$ , respectively. Find the asymptotic relative efficiency  $\sigma_2^2/\sigma_1^2$  of  $\tilde{\theta}_n$  with respect to  $\hat{\theta}_n$  and comment.

# Hypothesis Testing

## 6.1 Early History

Hypothesis testing at its early stage, from the 19th to the early 20th century, was concerned with hypotheses suggested by scientific theories in anthropometry, genetics, etc. The scientists in these disciplines wanted to evaluate the evidence provided by the data in support or against such hypotheses, called the *null hypotheses*.

The extent of departure from such a null hypothesis evidenced by the data was measured by a *test statistic* chosen on an ad hoc basis, which was considered *significant at level  $\alpha$*  if the tail probability “beyond” its observed value, under the null hypothesis, called the *p-value*, fell below  $\alpha$ . Whether to use the right tail or the left tail, or both tails for this purpose would be determined by the nature of the problem in a “supposedly obvious” manner. There was no formal basis for the choice of the test statistic or the direction of its tail probability, and the question of optimality of the commonly used test procedures was not addressed until [22] came up with the concept of an *alternative hypothesis  $H_1$*  against which the null hypothesis  $H_0$  was to be tested. This naturally led to the definitions of two types of error which are

- (i) deciding in favor of  $H_1$  when  $H_0$  is true (Type I error), and
- (ii) deciding in favor of  $H_0$  when  $H_1$  is true (Type II error).

Thus the problem of maximizing the

$$\begin{aligned}\text{Power} &= \text{The probability of correctly rejecting } H_0 \\ &= 1 - \text{the probability of Type II error}\end{aligned}$$

subject to the condition of the probability of Type I error not exceeding  $\alpha$  was defined and the solution to this problem was the *Neyman-Pearson Lemma*.

## 6.2 Basic Concepts

Let  $\{P_\theta, \theta \in \Theta\}$  be a family of probabilities on  $(\mathfrak{X}, \mathcal{A})$ , of which an unknown element  $P_\theta$  generates a random sample  $X$ . We consider two hypotheses,  $H_0: \theta \in \Theta_0$  and  $H_1: \theta \in \Theta_1$ , where  $\Theta_0$  and  $\Theta_1$  are disjoint subsets of  $\Theta$ . Based on the observation  $X$ , we want to take one

of two actions,  $a_0$ : Accept  $H_0$  (ie, decide  $\theta \in \Theta_0$ ) or  $a_1$ : Reject  $H_0$  and accept  $H_1$  (ie, decide  $\theta \in \Theta_1$ ). The problem of hypothesis testing was introduced in [Chapter 4](#). In this chapter we shall develop the methods of constructing optimal tests, mostly within the framework of exponential families, subject to the restrictions mentioned in [Section 4.6.2](#).

A nonrandomized decision rule is described by a function  $d: \mathfrak{X} \rightarrow \{a_0, a_1\}$ , or equivalently by  $C = \{x: d(x) = a_1\}$  which is called the *critical region* or *rejection region for  $H_0$* . The complement of  $C$  is  $A = \{x: d(x) = a_0\}$  which is called the *acceptance region for  $H_0$* .

A behavioral decision rule is described by a function  $\varphi: \mathfrak{X} \rightarrow [0, 1]$ , where for an  $x$ ,  $\varphi(x)$  is the probability of taking action  $a_1$  (ie, rejecting  $H_0$  when the observed value of  $X$  is  $x$ ). Such a function  $\varphi$  is called a *critical function*. A nonrandomized decision rule  $d$  with critical region  $C$  is equivalently described by a behavioral decision rule  $\varphi = I_C$ .

The critical region  $C$  of a nonrandomized decision rule and the sets  $\{x: \varphi(x) \leq c\}$  for  $c \in [0, 1]$  of the critical function  $\varphi$  of a behavioral decision rule must belong to  $\mathcal{A}$ .

We consider the  $0 - 1$  loss function, that is

$$L(\theta, a_0) = \begin{cases} 0 & \text{if } \theta \in \Theta_0 \\ 1 & \text{if } \theta \in \Theta_1 \end{cases}, \quad \text{and} \quad L(\theta, a_1) = \begin{cases} 1 & \text{if } \theta \in \Theta_0 \\ 0 & \text{if } \theta \in \Theta_1 \end{cases}.$$

Then the risk of a behavioral decision rule  $\varphi$  is

$$R(\theta, \varphi) = \begin{cases} E_\theta[\varphi(X)] := \text{Probability of Type I Error} & \text{if } \theta \in \Theta_0 \\ 1 - E_\theta[\varphi(X)] := \text{Probability of Type II Error} & \text{if } \theta \in \Theta_1 \end{cases}.$$

Thus the risk function is described in terms of the function

$$\beta_\varphi(\theta) = E_\theta[\varphi(X)].$$

On  $\Theta_0$ ,  $\beta_\varphi(\theta)$  is a measure of weakness of the test  $\varphi$  (probability of Type I error) and on  $\Theta_1$ ,  $\beta_\varphi(\theta)$  is a measure of strength of the test  $\varphi$  ( $1 - \beta_\varphi(\theta)$  – probability of Type II error) which is called the *power*.

In the theory of hypothesis testing, the hypotheses  $H_0: \theta \in \Theta_0$  and  $H_1: \theta \in \Theta_1$  are given asymmetric roles, because the Type I error receives more serious consideration than the Type II error. Optimization of  $\varphi$  is formulated as a problem of choosing  $\varphi$  so as to maximize  $\beta_\varphi(\theta)$  for  $\theta \in \Theta_1$  subject to the condition

$$\sup_{\theta \in \Theta_0} \beta_\varphi(\theta) \leq \alpha,$$

where  $\alpha \in [0, 1]$  is given. This  $\alpha$  is called the *level of significance* and  $\sup_{\theta \in \Theta_0} \beta_\varphi(\theta)$  is the *size* of the test  $\varphi$ .

We call " $H_0: \theta \in \Theta_0$ " the *null hypothesis* and " $H_1: \theta \in \Theta_1$ " the *alternative hypothesis*. If  $\Theta_0$  (and/or  $\Theta_1$ ) is a singleton set, then we call the null hypothesis  $H_0$  (and/or the alternative hypothesis  $H_1$ ) a *simple hypothesis*, otherwise  $H_0$  (and/or  $H_1$ ) is composite.

### 6.3 Simple Null Hypothesis vs Simple Alternative: Neyman-Pearson Lemma

Let  $P_0$  and  $P_1$  be two distinct probability distributions on  $(\mathfrak{X}, \mathcal{A})$  with pdf's/pmf's  $p_0$  and  $p_1$ , respectively. Let  $E_i$  denote  $E_{P_i}$ ,  $i = 0, 1$ .

**Definition 6.3.1.** A test  $\varphi$  is said to be *most powerful (MP)* at level  $\alpha$  for testing  $H_0: P = P_0$  against  $H_1: P = P_1$  if

- (a)  $E_0[\varphi(X)] \leq \alpha$ ; and
- (b) if any test  $\Psi$  satisfies (a), then  $E_1[\varphi(X)] \geq E_1[\Psi(X)]$ .

**Theorem 6.3.1** (Neyman-Pearson Lemma). *Suppose we are testing  $H_0: P = P_0$  against  $H_1: P = P_1$  and let  $0 \leq \alpha \leq 1$ .*

- (i) *Existence. There exists a test  $\varphi$  and a constant  $k$  such that  $E_0[\varphi(X)] = \alpha$  and*

$$\varphi(x) = \begin{cases} 1 & \text{if } p_1(x) > kp_0(x) \\ 0 & \text{if } p_1(x) < kp_0(x). \end{cases}$$

(ii) *Sufficiency. If a test  $\varphi$  satisfies the above conditions for some  $k$ , then it is MP at level  $\alpha$  for testing  $H_0: P = P_0$  against  $H_1: P = P_1$ .*

(iii) *Necessity. If  $\varphi$  is MP at level  $\alpha$  for testing  $H_0: P = P_0$  against  $H_1: P = P_1$ , then  $\varphi(x)$  is of the form given in (i) for some  $k$ .*

*Proof.*

- (i) Take  $0 < \alpha < 1$  since the proofs for  $\alpha = 0$  and  $1$  are straightforward. Let

$$F(c) = P_0[p_1(X) \leq cp_0(X)] = P_0[p_1(X)/p_0(X) \leq c]$$

be the cdf of the likelihood ratio  $T(X) = p_1(X)/p_0(X)$ , which is a bona fide rv under  $P_0$  because  $P_0[p_0(X) > 0] = 1$ . Since  $F$  is a cdf, there exists  $c_0$  such that

$$F(c_0 - 0) \leq 1 - \alpha \leq F(c_0).$$

Now define

$$\varphi(x) = \begin{cases} 1 & \text{if } p_1(x) > c_0 p_0(x) \\ \gamma & \text{if } p_1(x) = c_0 p_0(x) \\ 0 & \text{if } p_1(x) < c_0 p_0(x), \end{cases}$$

where  $\gamma$  can be assigned any value in  $[0, 1]$ , say  $\gamma = 1$ , if  $F$  is continuous at  $c_0$ , or else we take

$$\gamma = \frac{F(c_0) - (1 - \alpha)}{F(c_0) - F(c_0 - 0)}.$$

Thus part (i) of the N-P Lemma holds with  $k = c_0$  and  $0 \leq \gamma \leq 1$  as defined above, because in case of a jump in  $F(\cdot)$  at  $c_0$ ,

$$\begin{aligned} E_0[\varphi(X)] &= P_0[p_1(X) > c_0 p_0(X)] + \gamma P_0[p_1(X) = c_0 p_0(X)] \\ &= \{1 - F(c_0)\} + \gamma \{F(c_0) - F(c_0 - 0)\} = \alpha, \end{aligned}$$

and the same equality is obvious when  $\alpha$  is continuous at  $c_0$ .

- (ii) Let  $\varphi$  be a test satisfying the conditions in (i) and let  $\psi$  be another level  $\alpha$  test for  $H_0$  (ie,  $E_0[\psi(X)] \leq \alpha$ ). Let

$$S^+ = \{x: \varphi(x) - \psi(x) > 0\} \text{ and } S^- = \{x: \varphi(x) - \psi(x) < 0\}.$$

Then for  $x \in S^+$ ,  $\varphi(x) > 0$ , so  $p_1(x) \geq kp_0(x)$  and for  $x \in S^-$ ,  $\varphi(x) < 1$ , so  $p_1(x) \leq kp_0(x)$ . Thus

$$\begin{aligned} \{\varphi(x) - \psi(x)\}\{p_1(x) - kp_0(x)\} &\geq 0 \text{ for all } x \in S^+ \cup S^- \text{ and} \\ &= 0 \text{ for all other } x. \end{aligned}$$

Hence

$$\begin{aligned} E_1[\varphi(X)] - E_1[\psi(X)] &= \int \{\varphi(x) - \psi(x)\}p_1(x) dx \\ &= \int \{\varphi(x) - \psi(x)\}\{p_1(x) - kp_0(x)\} dx \\ &\quad + k \int \{\varphi(x) - \psi(x)\}p_0(x) dx \\ &\geq 0 + k\{E_0[\varphi(X)] - E_0[\psi(X)]\} \geq 0. \end{aligned}$$

- (iii) Suppose that  $\psi$  is MP at level  $\alpha$  for testing  $H_0: P = P_0$  vs  $H_1: P = P_1$ , and let  $\varphi$  be a test satisfying the conditions in (i) with  $k = k^*$ . Define  $S^+$  and  $S^-$  as in the proof of part (ii). Then  $\psi(x) \neq \varphi(x)$  for  $x \in S^+ \cup S^-$ . Let

$$\begin{aligned} S &= (S^+ \cup S^-) \cap \{x: p_1(x) \neq k^*p_0(x)\} \\ &= \{x: \psi(x) \neq \varphi(x) \text{ and } p_1(x) \neq k^*p_0(x)\}. \end{aligned}$$

Then  $P_0(S) = P_1(S) = 0$  implies  $\psi(x) = \varphi(x)$  w.p. 1 under  $P_0$  and  $P_1$  whenever  $p_1(x) \neq k^*p_0(x)$ , showing that  $\psi$  satisfies condition (i) for  $k = k^*$  w.p. 1 under  $P_0$  and  $P_1$ , proving part (iii) of the lemma. To prove  $P_0(S) = P_1(S) = 0$ , it is enough to show that  $\int_S dx = 0$ . Since  $(\varphi - \psi)(p_1 - k^*p_0) > 0$  on  $S$  and = 0 on  $S^c$ ,

$$\begin{aligned} \int (\varphi - \psi)(p_1 - k^*p_0) &= \int_S (\varphi - \psi)(p_1 - k^*p_0) > 0, \text{ ie,} \\ \int (\varphi - \psi)p_1 &> k^* \int (\varphi - \psi)p_0, \end{aligned}$$

implying  $E_1[\varphi(X)] > E_1[\psi(X)]$ , which is a contradiction.

□

### Corollary to the N-P Lemma

Let  $\beta = E_1[\varphi(X)]$  where  $\varphi$  is the MP test for  $H_0: P = P_0$  vs  $H_1: P = P_1$  at level  $\alpha < 1$ . Then  $\beta > \alpha$ .

*Proof.* Consider the test  $\varphi_0(x) = \alpha$  for all  $x$ . Then  $E_0[\varphi_0(X)] = E_1[\varphi_0(X)] = \alpha$ . Since  $\varphi_0$  is a level  $\alpha$  test which does not satisfy the necessary condition for an MP test (given by part

(iii) of the N-P Lemma) it is *not* an MP level  $\alpha$  test, whereas  $\varphi$  is. Hence

$$\beta = E_1[\varphi(X)] > E_1[\varphi_0(X)] = \alpha.$$

□

**Note.**  $\varphi_0(x) = \alpha$  for all  $x$  may satisfy the necessary condition if  $p_1(x) \equiv kp_0(x)$  for some  $k$ ; but then  $k = 1$ , contradicting  $P_0 \neq P_1$ .

*Remark 6.3.1.* For two probabilities  $P_0, P_1$  on  $(\mathfrak{X}, \mathcal{A})$ , we call  $p_1(x)/p_0(x)$  the *likelihood ratio* of  $P_1$  to  $P_0$  at  $x$ . The N-P Lemma expresses the MP level  $\alpha$  test for  $H_0: P = P_0$  vs  $H_1: P = P_1$  in terms of the likelihood ratio, or equivalently, in terms of the log likelihood ratio  $L(x) = \log(p_1(x)/p_0(x))$

$$\phi(x) = \begin{cases} 0 & \text{if } L(x) < k \\ \gamma & \text{if } L(x) = k \\ 1 & \text{if } L(x) > k, \end{cases}$$

where  $k$  and  $0 \leq \gamma \leq 1$  are determined by the condition  $E_0[\phi(X)] = \alpha$ .

Although the N-P Lemma can be used to construct MP level  $\alpha$  tests for arbitrary  $P_0$  vs  $P_1$  such as testing whether the data came from  $N(0, 1)$  or from  $Cauchy(0, 1)$ , our interest mostly lies in testing for one value of the parameter against another *within a parametric family*. We illustrate the use of the N-P Lemma with two examples, namely, testing for the mean  $\theta$  of  $Poi(\theta)$  and the mean  $\theta$  of  $N(\theta, 1)$ .

**Example 6.3.1.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $Poisson(\theta)$ , where  $\theta > 0$  is unknown. We want the MP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1$ .

*Solution.* Here the log likelihood ratio is

$$\log \frac{p_{\theta_1}(\mathbf{x})}{p_{\theta_0}(\mathbf{x})} = -n(\theta_1 - \theta_0) + (\log \theta_1 - \log \theta_0)t_n, \text{ where } t_n = \sum_{i=1}^n x_i.$$

**(i)** Hence the MP level  $\alpha$  test for  $H_0$  vs  $H_1$  is

$$\varphi_1(\mathbf{x}) = \begin{cases} 0 & \text{if } t_n < c_1 \\ \gamma_1 & \text{if } t_n = c_1 \\ 1 & \text{if } t_n > c_1, \end{cases}$$

where  $c_1$  and  $0 \leq \gamma_1 \leq 1$  are determined by  $P_{\theta_0}[T_n > c_1] + \gamma_1 P_{\theta_0}[T_n = c_1] = \alpha$  if  $\theta_1 > \theta_0$  and

$$\varphi_2(\mathbf{x}) = \begin{cases} 0 & \text{if } t_n > c_2 \\ \gamma_2 & \text{if } t_n = c_2 \\ 1 & \text{if } t_n < c_2, \end{cases}$$

where  $c_2$  and  $0 \leq \gamma_2 \leq 1$  are determined by  $P_{\theta_0}[T_n < c_2] + \gamma_2 P_{\theta_0}[T_n = c_2] = \alpha$ , if  $\theta_1 < \theta_0$ .

**(ii)** We illustrate the determination of  $c_1$  and  $\gamma_1$  in the first case by considering  $H_0: \theta = \theta_0 = 0.5$  and  $H_1: \theta = \theta_1 = 1.0$  for  $n = 5$  and  $\alpha = 0.5$ . Under  $P_{\theta_0}$ ,  $T = \sum_{i=1}^5 X_i \sim Poisson(2.5)$  (ie,  $p_{\theta_0}(t) = e^{-2.5}(2.5)^t/t!$ ,  $t = 0, 1, 2, \dots$ ). Calculating these probabilities, we have  $P_{\theta_0}[T > 5] = 0.0420$  and  $P_{\theta_0}[T = 5] = 0.0668$  and then solve the equation:

$$P_{\theta_0}[T > 5] + \gamma P_{\theta_0}[T = 5] = 0.05.$$

Thus  $c_1 = 5$  and  $\gamma_1 = 0.12$ , so the MP test for  $H_0: \theta = 0.5$  vs  $H_1: \theta = 1.0$  at level  $\alpha = 0.05$  is

$$\varphi_1(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_1^5 x_i < 5 \\ 0.12 & \text{if } \sum_1^5 x_i = 5 \\ 1 & \text{if } \sum_1^5 x_i > 5 \end{cases}$$

(iii) The power function of  $\varphi_1$  is

$$\beta_1(\theta) = P_\theta[T > 5] + 0.12P_\theta[T = 5] = \sum_{t=6}^{\infty} e^{-\lambda} \lambda^t / t! + 0.12e^{-\lambda} \lambda^5 / 5!, \text{ where}$$

$$\lambda = 5\theta \text{ and } \beta'_1(\theta) = 5e^{-\lambda} \left[ 0.12 \frac{\lambda^4}{4!} + 0.88 \frac{\lambda^5}{5!} \right] > 0.$$

(iv) The MP level  $\alpha$  test  $\varphi_1$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1 (> \theta_0)$  does not depend on  $\theta_1$ , so long as  $\theta_1 > \theta_0$ . Therefore, the MP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1 (> \theta_0)$  is the *Uniformly Most Powerful (UMP)* level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1^*: \theta > \theta_0$ . Again, since  $\beta'_1(\theta) > 0$ ,  $\beta_1(\theta_0) = \alpha$  implies  $\beta_1(\theta) < \alpha$  for all  $\theta < \theta_0$ . Thus the UMP level  $\alpha$  test  $\varphi_1$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta > \theta_0$  is the UMP level  $\alpha$  test for the composite null hypothesis  $H_0^*: \theta \leq \theta_0$  vs the composite alternative  $H_1^*: \theta > \theta_0$ .

In the same way  $\varphi_2$  can be shown to be a UMP level  $\alpha$  test for  $H_0^*: \theta \geq \theta_0$  vs  $H_1^*: \theta < \theta_0$ .

**Example 6.3.2.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $N(\theta, \sigma^2)$  where  $\theta$  is unknown but  $\sigma^2 > 0$  is known. We want the MP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1$ .

*Solution.* Here the log likelihood ratio is

$$\log(p_{\theta_1, \sigma}(\mathbf{x})/p_{\theta_0, \sigma}(\mathbf{x})) = \frac{n(\theta_1 - \theta_0)}{\sigma^2} \bar{x}_n - \frac{n(\theta_1^2 - \theta_0^2)}{\sigma^2}, \quad \text{where } \bar{x}_n = n^{-1} \sum_{i=1}^n x_i.$$

Hence the MP level  $\alpha$  test for  $H_0$  vs  $H_1$  is

$$\varphi_1(\mathbf{x}) = \begin{cases} 0 & \bar{x}_n < c_1, \text{ if } \theta_1 > \theta_0, \text{ and} \\ 1 & \bar{x}_n \geq c_1 \end{cases}$$

$$\varphi_2(\mathbf{x}) = \begin{cases} 0 & \bar{x}_n > c_2, \text{ if } \theta_1 < \theta_0, \\ 1 & \bar{x}_n \leq c_2 \end{cases}$$

where  $c_1$  is determined by  $P_{\theta_0}[\bar{X}_n \geq c_1] = \alpha$  and  $c_2$  is determined by  $P_{\theta_0}[\bar{X}_n \leq c_2] = \alpha$ . The constants  $c_1$  and  $c_2$  are

$$c_1 = \theta_0 + (\sigma/\sqrt{n})\Phi^{-1}(1 - \alpha) \text{ and } c_2 = \theta_0 - (\sigma/\sqrt{n})\Phi^{-1}(1 - \alpha),$$

where  $\Phi$  is the cdf of  $N(0, 1)$ .

The power function of  $\varphi_1$  is

$$\begin{aligned}\beta_1(\theta) &= P_\theta[\bar{X}_n \geq \theta_0 + (\sigma/\sqrt{n})\Phi^{-1}(1-\alpha)] = \Phi(\Phi^{-1}(\alpha) + \sqrt{n}(\theta - \theta_0)/\sigma), \text{ and} \\ \beta'_1(\theta) &= (\sqrt{n}/\sigma)\Phi'(\Phi^{-1}(\alpha) + \sqrt{n}(\theta - \theta_0)/\sigma) > 0.\end{aligned}$$

Using the same argument as in [Example 6.3.1](#), we see that the MP level  $\alpha$  test  $\varphi_1$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1 (> \theta_0)$  is also the UMP level  $\alpha$  test for the composite null hypothesis  $H_0^*: \theta \leq \theta_0$  vs the composite alternative  $H_1^*: \theta > \theta_0$ .

The MP level  $\alpha$  test  $\varphi_2$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1$  if  $\theta_1 < \theta_0$ , in the same way.

## 6.4 UMP Tests for One-Sided Hypotheses Against One-Sided Alternatives in Monotone Likelihood Ratio Families

We begin this section with the definition of UMP tests.

**Definition 6.4.1.** A test  $\varphi$  is a UMP test at level  $\alpha$  for  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$  if

- (i)  $\sup_{\theta \in \Theta_0} E_\theta[\varphi(X)] \leq \alpha$ , and
- (ii)  $E_\theta[\varphi(X)] \geq E_\theta[\psi(X)]$  for all  $\theta \in \Theta_1$ , whenever  $\psi$  satisfies (i).

Although requirement (ii) is very stringent, UMP tests do exist in certain types of situations. In the previous section we have seen that for  $Poisson(\theta)$ , the MP test at level  $\alpha$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1 (> \theta_0)$  has the stronger property of being the UMP level  $\alpha$  test for the composite hypothesis  $H_0: \theta \leq \theta_0$  vs the composite alternative  $H_1: \theta > \theta_0$ . This actually holds in a much wider context.

**Definition 6.4.2.** Let  $\{P_\theta, \theta \in \Theta\}$  be a family of probabilities on  $(\mathfrak{X}, \mathcal{A})$  and let  $p_\theta$  denote the pdf or pmf corresponding to  $P_\theta$  where  $\Theta = \mathbb{R}$  or  $\Theta$  is an interval in  $\mathbb{R}$ . Such a family  $\{p_\theta\}$  is said to be a monotone likelihood ratio (MLR) family if there exists a real-valued statistic  $T(x)$  such that for any  $\theta_1 < \theta_2$  in  $\Theta$ ,  $p_{\theta_2}(x)/p_{\theta_1}(x)$  is a nondecreasing function of  $T(x)$ . [If  $p_{\theta_1}(x) = 0 < p_{\theta_2}(x)$ , define  $p_{\theta_2}(x)/p_{\theta_1}(x) = +\infty$ .]

**Example 6.4.1.** Let  $p_\theta(x) = c(\theta) \exp[Q(\theta)T(x)]h(x)$  where  $Q(\theta)$  is a nondecreasing function. Then  $\{p_\theta\}$  is an MLR family. This includes

- (a)  $p_\theta(x) = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_1^n (x_i - \theta)^2\right], x_i \in \mathbb{R}, \theta \in \mathbb{R}, \sigma^2 \text{ fixed.}$
- (b)  $p_\theta(x) = (2\pi\theta)^{-n/2} \exp\left[-\frac{1}{2\theta} \sum_1^n (x_i - \mu)^2\right], x_i \in \mathbb{R}, \theta \in \mathbb{R}^+, \mu \text{ fixed.}$
- (c)  $p_\theta(x) = e^{-n\theta} \theta^{\sum_1^n x_i} / \prod_1^n x_i!, x_i \in \{0, 1, 2, \dots\}, \theta \in \mathbb{R}^+.$
- (d)  $p_\theta(x) = \theta^{\sum_1^n x_i} (1 - \theta)^{n - \sum_1^n x_i}, x_i \in \{0, 1\}, \theta \in (0, 1).$

**Example 6.4.2.** The hypergeometric distribution  $\mathcal{H}(n, N, \theta)$  with  $p_\theta(x) = \binom{\theta}{x} \binom{N-\theta}{n-x} / \binom{N}{n}$ ,  $x = \max(0, \theta + n - N), \dots, \min(n, \theta)$ .

**Example 6.4.3.** The family of Cauchy distributions  $\mathcal{C}(\theta, 1)$  with  $p_\theta(x) = \frac{1}{\pi[1+(x-\theta)^2]}$ ,  $x \in \mathbb{R}, \theta \in \mathbb{R}$  is *not* an MLR family.

**Theorem 6.4.1.** Suppose that the family of pdf's or pmf's  $\{p_\theta, \theta \in \mathbb{R}\}$  has MLR property in  $T(x)$ . Then

- (i) There exists a UMP level  $\alpha$  test for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  given by

$$\varphi(x) = \begin{cases} 0 & \text{if } T(x) < c \\ \gamma & \text{if } T(x) = c, \text{ where } c \text{ and } 0 \leq \gamma \leq 1 \text{ are determined by } E_{\theta_0}[\varphi(X)] = \alpha. \\ 1 & \text{if } T(x) > c \end{cases}$$

- (ii) The power function  $\beta(\theta) = E_\theta[\varphi(X)]$  is strictly increasing at all  $\theta$  for which  $\beta(\theta) < 1$ .

- (iii) For all  $\theta'$ , this  $\varphi$  is UMP at level  $\alpha' = \beta(\theta')$  for  $H'_0: \theta \leq \theta'$  vs  $H'_1: \theta > \theta'$ .

- (iv) This test minimizes  $E_\theta[\psi(X)]$  at all  $\theta < \theta_0$  among all tests  $\psi$  for which  $E_{\theta_0}[\psi(X)] = \alpha$ .

*Proof.* Parts (i) and (ii). First consider the simple vs simple case with  $H_0^*: \theta = \theta_0$  vs  $H_1^*: \theta = \theta_1$  where  $\theta_1 > \theta_0$  is fixed. Then by the N-P Lemma, the MP level  $\alpha$  test for  $H_0^*$  vs  $H_1^*$  rejects  $H_0$  for large values of  $p_{\theta_1}(x)/p_{\theta_0}(x)$  (ie, for large values of  $T(x)$ ) by the MLR property. Moreover, by the existence part of the N-P Lemma, there exist  $c$  and  $0 \leq \gamma \leq 1$  such that the test

$$\varphi(x) = I_{(c, \infty)}[T(X)] + \gamma I_{\{c\}}[T(X)] \text{ satisfies } E_{\theta_0}[\varphi(X)] = \alpha.$$

Since the forms:

$$\varphi(x) = \begin{cases} 0, & T(x) < c \\ 1, & T(x) > c \end{cases} \text{ and } \varphi(x) = \begin{cases} 0, & p_{\theta''}(x) < kp_{\theta'}(x) \\ 1, & p_{\theta''}(x) > kp_{\theta'}(x) \end{cases}$$

are equivalent for any  $\theta' < \theta''$ , this test is UMP at level  $\alpha' = \beta(\theta')$  for testing  $H_0^{**}: \theta = \theta'$  vs  $H_1^{**}: \theta = \theta''$  whenever  $\theta' < \theta''$  (by the sufficiency part of the N-P Lemma). Next, note that by the corollary to the N-P Lemma,  $\beta(\theta'') > \alpha' = \beta(\theta')$  if  $\alpha' < 1$ , which proves that  $\beta(\theta)$  is strictly increasing, so long as it is  $< 1$ . This proves part (ii).

Now note that for this test,  $\beta(\theta) = E_\theta[\varphi(X)] \leq \alpha$  for all  $\theta \leq \theta_0$ , which makes  $\varphi$  a level  $\alpha$  test for  $H_0: \theta \leq \theta_0$ . Let

$$\begin{aligned} \Psi_\alpha &= \left\{ \text{all tests } \psi \text{ such that } \sup_{\theta \leq \theta_0} E_\theta[\psi(X)] \leq \alpha \right\} \text{ and} \\ \Psi_\alpha^* &= \{ \text{all tests } \psi \text{ such that } E_{\theta_0}[\psi(X)] \leq \alpha \}. \end{aligned}$$

Then  $\Psi_\alpha \subset \Psi_\alpha^*$ . We have shown that

$$\varphi \in \Psi_\alpha \text{ and } E_{\theta_1}[\varphi(X)] \geq E_{\theta_1}[\psi(X)] \text{ for all } \psi \in \Psi_\alpha^*.$$

Hence  $E_{\theta_1}[\varphi(X)] \geq E_{\theta_1}[\psi(X)]$  for all  $\psi \in \Psi_\alpha$ .

This makes  $\varphi$  the MP test at level  $\alpha$  for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$ . Finally, since  $\varphi$  does not depend on  $\theta_1 > \theta_0$ , it is UMP at level  $\alpha$  for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$ .

The proof of part (iii) is analogous, and part (iv) is proved by observing that power is minimized if all inequalities are reversed in the N-P Lemma.  $\square$

## 6.5 Unbiased Tests

A behavioral test with critical function  $\varphi$  is said to be an unbiased test at level  $\alpha$  for  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$  if (i)  $E_\theta[\varphi(X)] \leq \alpha$  for all  $\theta \in \Theta_0$  and (ii)  $E_\theta[\varphi(X)] \geq \alpha$  for all  $\theta \in \Theta_1$ . An unbiased test rejects the null hypothesis with at least as much probability when it is false as when it is true.

**Definition 6.5.1.** A test  $\varphi$  is said to be UMP unbiased test at level  $\alpha$  for  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$  if  $\varphi$  is an unbiased test at level  $\alpha$  for  $H_0$  vs  $H_1$ , and if  $E_\theta[\varphi(X)] \geq E_\theta[\psi(X)]$  for all  $\theta \in \Theta_1$  whenever  $\psi$  is also an unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$ .

For a large class of problems, a UMP test does not exist, but a UMP unbiased test does exist.

**Example 6.5.1.** Let  $X \sim \mathcal{N}(\theta, 1)$ ,  $H_0: \theta = \theta_0$  and  $H_1: \theta \neq \theta_0$ . Here the tests

$$\varphi_1(x) = I_{[\theta_0 + \Phi^{-1}(1-\alpha), \infty)}(x) \text{ and } \varphi_2(x) = I_{(-\infty, \theta_0 - \Phi^{-1}(1-\alpha))}(x),$$

where  $\Phi$  is the cdf of  $\mathcal{N}(0, 1)$ , are, respectively, the UMP level  $\alpha$  tests for  $H_0: \theta = \theta_0$  vs  $H_1^+: \theta > \theta_0$  and  $H_1^-: \theta < \theta_0$ . Moreover, it follows from [Theorem 6.4.1](#), part (ii) that  $E_\theta[\varphi_1(x)] < \alpha$  for all  $\theta < \theta_0$  and  $E_\theta[\varphi_2(x)] < \alpha$  for all  $\theta > \theta_0$  (because the family  $\mathcal{N}(\theta, 1)$  has the MLR property), which shows that a UMP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  does not exist.

For testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  at level  $\alpha$ , an unbiased test  $\varphi$  must satisfy:  $E_{\theta_0}[\varphi(X)] \leq \alpha$  and  $E_\theta[\varphi(X)] \geq \alpha$  for all  $\theta \neq \theta_0$ , so neither  $\varphi_1$  nor  $\varphi_2$  in the above example is an unbiased test. If we restrict to the class of unbiased level  $\alpha$  tests then  $\varphi_1, \varphi_2$  would not qualify, but in the restricted class a UMP test does exist.

Suppose that the power functions  $\beta_\varphi(\theta) = E_\theta[\varphi(X)]$  of all tests  $\varphi$  are differentiable at  $\theta_0$  (as in the case of exponential families). Then an unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  must satisfy

$$E_{\theta_0}[\varphi(X)] = \alpha \quad \text{and} \quad \beta'_\varphi(\theta_0) = 0. \quad (1)$$

Let  $\mathcal{C}_0$  denote the class of unbiased level  $\alpha$  tests for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  and let  $\mathcal{C}_1$  denote the class of all tests satisfying Eq. (1). Then  $\mathcal{C}_0 \subset \mathcal{C}_1$  and therefore, if a test  $\varphi_0 \in \mathcal{C}_0$  is UMP among all tests in  $\mathcal{C}_1$ , then it is a UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$ .

We now examine the case of a single-parameter exponential family  $\{P_\theta, \theta \in \Theta \subset \mathbb{R}\}$  with pdf/pmf given by

$$p_X(x; \theta) = C_0(\theta) \exp[\theta T(x)] h_0(x), \quad x \in \mathfrak{X}, \quad (2)$$

and  $\mathbf{X} = (X_1, \dots, X_n)$  being a random sample from  $P_\theta$ . Then  $T = \sum_{i=1}^n T(X_i)$  is sufficient for  $\theta$  in  $(X_1, \dots, X_n)$  and is distributed with pdf/pmf

$$p_T(t; \theta) = C(\theta) \exp[\theta t] h(t), \quad t \in \mathcal{T}. \quad (3)$$

We know that for all tests  $\varphi$ , the power functions

$$\beta_\varphi(\theta) = C(\theta) \int \varphi(t) \exp[\theta t] h(t) dt \quad (4)$$

are differentiable at all interior points of the natural parameter space, which is an interval, and if  $\theta_0$  is such a point, then the differentiation can be carried out under the integral, that is

$$\begin{aligned}\beta'_\varphi(\theta_0) &= C'(\theta_0) \int \varphi(t) \exp[\theta_0 t] h(t) dt + C(\theta_0) \int t \varphi(t) \exp[\theta_0 t] h(t) dt \\ &= \frac{C'(\theta_0)}{C(\theta_0)} E_{\theta_0}[\varphi(T)] + E_{\theta_0}[T \varphi(T)].\end{aligned}$$

Moreover, for  $\varphi(t) \equiv \alpha$ , we have  $\beta_\varphi(\theta) = \alpha$  for all  $\theta$ , and therefore,  $\beta'_\varphi(\theta_0) = 0$ . Thus the above expression becomes:  $0 = \frac{C'(\theta_0)}{C(\theta_0)} \alpha + \alpha E_{\theta_0}[T]$ . Hence  $C'(\theta_0)/C(\theta_0) = -E_{\theta_0}[T]$ , and the formula for  $\beta'_\varphi(\theta_0)$  for an arbitrary level  $\alpha$  test becomes:

$$\beta'_\varphi(\theta_0) = -\alpha E_{\theta_0}[T] + E_{\theta_0}[T \varphi(T)]. \quad (5)$$

The problem of finding the UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  in a single-parameter exponential family, where  $\theta_0$  is an interior point of the natural parameter space, now leads to the following restatement of Eq. (1), using the formula (5) for  $\beta'_\varphi(\theta_0)$ .

Among all tests  $\varphi \in \mathcal{C}_1$  satisfying

$$E_{\theta_0}[\varphi(T)] = \int \varphi(t) f_1(t) dt = \alpha, \text{ and} \quad (6)$$

$$E_{\theta_0}[T \varphi(T)] = \int \varphi(t) f_2(t) dt = \alpha E_{\theta_0}[T] \quad (7)$$

with

$$f_1(t) = C(\theta_0) \exp[\theta_0 t] h(t), \quad (8)$$

$$f_2(t) = C(\theta_0) t \exp[\theta_0 t] h(t) \quad (9)$$

find  $\varphi_0$  which maximizes  $E_{\theta_1}[\varphi(T)] = \int \varphi(t) f_3(t) dt$  with

$$f_3(t) = C(\theta_1) \exp[\theta_1 t] h(t) \quad (10)$$

for a fixed  $\theta_1 \neq \theta_0$ .

If this  $\varphi_0$  is in the smaller class  $\mathcal{C}_0$  of unbiased level  $\alpha$  tests for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  and if  $\varphi_0$  does not depend on the specific  $\theta_1 \neq \theta_0$  used in the above optimization problem, then  $\varphi_0$  is a UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$ .

We next look at the problem of testing  $H_0: \theta_1 \leq \theta \leq \theta_2$  vs  $H_1: \theta \notin [\theta_1, \theta_2]$  in an exponential family, where  $\theta_1, \theta_2$  are interior points of the natural parameter space. More generally, first consider the problem of testing  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$ , where  $\Theta_0$  and  $\Theta_1$  have a common boundary  $\omega$ . Since the power functions  $\beta_\varphi(\theta) = E_\theta[\varphi(X)]$  of all tests are continuous in  $\theta$ ,  $E_\theta[\varphi(X)] = \alpha$  for all  $\theta \in \omega$  must hold for all unbiased level  $\alpha$  tests for  $H_0$  vs  $H_1$ .

**Definition 6.5.2.** A test  $\varphi$  satisfying  $E_\theta[\varphi(X)] = \alpha$  for all  $\theta \in \omega$  is said to be *similar* of size  $\alpha$  on  $\omega$ .

Since the class of similar tests of size  $\alpha$  on  $\omega$  includes the class of unbiased level  $\alpha$  tests for  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$ , the following holds.

**Lemma 6.5.1.** *If the power functions of all tests are continuous in  $\theta$  and if a test  $\varphi_0$  is UMP among all similar tests of size  $\alpha$  on the common boundary  $\omega$  of  $\Theta_0$  and  $\Theta_1$ , then  $\varphi_0$  is UMP unbiased level  $\alpha$  test for  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \Theta_1$  provided that  $\varphi_0$  is a level  $\alpha$  test.*

*Proof.* Clearly,  $\varphi_0$  is UMP on  $\Theta_1$  among all unbiased level  $\alpha$  tests. Therefore, we only need to verify that  $\varphi_0$  is an unbiased level  $\alpha$  test. For this, we compare  $\varphi_0$  with  $\varphi^*(x) \equiv \alpha$  (which is a similar test of size  $\alpha$ ), to see that  $E_\theta[\varphi_0(X)] \geq E_\theta[\varphi^*(X)] = \alpha$  for all  $\theta \in \Theta_1$ .  $\square$

We now consider the problem of testing  $H_0: \theta_1 \leq \theta \leq \theta_2$  vs  $H_1: \theta \notin [\theta_1, \theta_2]$  in an exponential family  $\{P_\theta, \theta \in \Theta \subset \mathbb{R}\}$  with pdf/pmf  $p_X(x; \theta)$  given by Eq. (2), based on a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  in which  $T = \sum_{i=1}^n T(X_i)$  is sufficient for  $\theta$  with pdf/pmf given  $f_T(t; \theta)$  given by Eq. (3). We also know that for all tests  $\varphi$ , the power function  $\beta_\varphi(\theta)$  given by Eq. (4) is continuous in  $\theta$ .

To find a UMP unbiased level  $\alpha$  test for  $H_0: \theta \in \Theta_0 = [\theta_1, \theta_2]$  vs  $H_1: \theta \in \Theta_1 = [\theta_1, \theta_2]^c$ , we note that  $\Theta_0$  and  $\Theta_1$  have a common boundary  $\omega = \{\theta_1, \theta_2\}$ , and all tests have continuous power functions. So the above Lemma is applicable and we only need to look among tests which are similar of size  $\alpha$  on  $\omega = \{\theta_1, \theta_2\}$  (ie, among tests satisfying  $E_{\theta_i}[\varphi(T)] = \alpha$ ,  $i = 1, 2$ ). This leads to the problem of maximizing  $\int \varphi(t)f_3(t) dt$  subject to the conditions  $\int \varphi(t)f_i(t) dt = \alpha$ ,  $i = 1, 2$ , where

$$f_i(t) = C(\theta_i) \exp[\theta_i t] h(t), \quad i = 1, 2, \text{ and} \quad (11)$$

$$f_3(t) = C(\theta) \exp[\theta t] h(t) \quad \text{for } \theta \notin [\theta_1, \theta_2]. \quad (12)$$

Thus, our search for UMP unbiased level  $\alpha$  tests for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  and  $H_0: \theta_1 \leq \theta \leq \theta_2$  vs  $H_1: \theta \notin [\theta_1, \theta_2]$  in a single-parameter exponential family, both lead to the problem of maximizing  $\int \varphi(t)f_3(t) dt$  with respect to  $\varphi$  subject to the conditions  $\int \varphi(t)f_i(t) dt = c_i$ ,  $i = 1, 2$ , where the constants are  $c_1 = \alpha$ ,  $c_2 = \alpha E_{\theta_0}[T]$  and  $f_1, f_2, f_3$  given by Eqs. (8), (9), and (10) in the first problem and  $c_1 = c_2 = \alpha$  and  $f_1, f_2, f_3$  given by Eqs. (11) and (12) in the second problem. This calls for a generalization of the Neyman-Pearson Lemma [23].

## 6.6 Generalized Neyman-Pearson Lemma

**Theorem 6.6.1.** *Let  $f_1, \dots, f_m, f_{m+1}$  be real-valued functions defined on an Euclidean space ( $\mathfrak{X} = \mathbb{R}^n, \mathcal{A}$ ) for which  $\int_{\mathfrak{X}} f_i(x) dx$  (or  $\sum_{x \in \mathfrak{X}} f_i(x)$  in the discrete case) exist and are finite for  $i = 1, \dots, m+1$ . Let  $\mathcal{F}$  denote the class of all critical functions (ie  $\varphi: \mathfrak{X} \rightarrow [0, 1]$ ). Suppose that for given constants  $c_1, \dots, c_m$ , there exists  $\varphi \in \mathcal{F}$  such that  $\int \varphi f_i dx = c_i$ ,  $i = 1, \dots, m$ . Let*

$$\mathcal{C} = \left\{ \varphi \in \mathcal{F}: \int \varphi f_i dx = c_i, \quad i = 1, \dots, m \right\}$$

which is nonempty.

(i) If there exists  $\varphi \in \mathcal{C}$  such that for some constants  $k_1, \dots, k_m$ ,

$$\varphi(x) = \begin{cases} 0 & \text{iff } m+1(x) < \sum_{i=1}^m k_i f_i(x) \\ 1 & \text{iff } m+1(x) > \sum_{i=1}^m k_i f_i(x), \end{cases}$$

then  $\int \varphi f_{m+1} dx \geq \psi f_{m+1} dx$  for all  $\psi \in \mathcal{C}$ . This provides a sufficient condition for maximization of  $\int \varphi f_{m+1} dx$  in  $\mathcal{C}$ .

- (ii) If there exists  $\varphi \in \mathcal{C}$  such that for constants  $k_1, \dots, k_m \geq 0$ ,  $\varphi(x)$  is of the form given in (i), then  $\int \varphi f_{m+1} dx \geq \int \psi f_{m+1} dx$  for all  $\psi \in \mathcal{C}' = \{\psi \in \mathcal{F}: \int \varphi f_i dx \leq c_i, i = 1, \dots, m\}$ . This is an extension of (i).

*Proof.* Let  $\varphi \in \mathcal{C}$  be of the form given in (i) and  $\psi \in \mathcal{F}$ . Let

$$S^+ = \left\{ x: f_{m+1}(x) > \sum_{i=1}^m k_i f_i(x) \right\} \text{ and } S^- = \left\{ x: f_{m+1}(x) < \sum_{i=1}^m k_i f_i(x) \right\}.$$

For  $x \in S^+$ ,  $\varphi(x) - \psi(x) = 1 - \psi(x) \geq 0$  and  $f_{m+1} - \sum_{i=1}^m k_i f_i(x) > 0$  and for  $x \in S^-$ , both of these inequalities are reversed. Hence

$$\{\varphi(x) - \psi(x)\} \left\{ f_{m+1}(x) - \sum_{i=1}^m k_i f_i(x) \right\} \geq 0 \quad \text{for all } x \in S^+ \cup S^-.$$

Thus,

$$\begin{aligned} \int_{S^+ \cup S^-} (\varphi - \psi) \left( f_{m+1} - \sum_{i=1}^m k_i f_i \right) dx &\geq 0, \text{ ie,} \\ \int (\varphi - \psi) f_{m+1} dx &\geq \sum_{i=1}^m k_i \int (\varphi - \psi) f_i dx. \end{aligned}$$

If  $\psi \in \mathcal{C}$ , then  $\sum_{i=1}^m k_i \int (\varphi - \psi) f_i dx = 0$ ; if  $\psi \in \mathcal{C}'$  and  $k_1, \dots, k_m \geq 0$ , then  $\sum_{i=1}^m k_i \int (\varphi - \psi) f_i dx \geq 0$ . Hence

$$\int \varphi f_{m+1} dx \geq \int \psi f_{m+1} dx \text{ if } \psi \in \mathcal{C}, \text{ or if } \psi \in \mathcal{C}' \text{ and } k_1, \dots, k_m \geq 0.$$

□

## 6.7 UMP Unbiased Tests for Two-Sided Problems

We now apply the Generalized N-P Lemma to find unbiased level  $\alpha$  tests for the two problems under discussion.

### 6.7.1 UMP Unbiased Test for $H_0: \theta = \theta_0$ vs $H_1: \theta \neq \theta_0$ in a Single-Parameter Exponential Family

By the Generalized N-P Lemma, the UMP unbiased level  $\alpha$  test for this problem is given by:

$$\varphi_0(t) = \begin{cases} 0 & \text{if } f_3(t) < k_1 f_1(t) + k_2 f_2(t) \\ \gamma(t) & \text{if } f_3(t) = k_1 f_1(t) + k_2 f_2(t) \\ 1 & \text{if } f_3(t) > k_1 f_1(t) + k_2 f_2(t), \end{cases}$$

where  $f_1, f_2$ , and  $f_3$  are given by Eqs. (8), (9), and (10) and  $k_1, k_2$ , and  $0 \leq \gamma(t) \leq 1$  are such that

$$E_{\theta_0}[\varphi_0(T)] = \alpha \quad \text{and} \quad E_{\theta_0}[T\varphi_0(T)] = \alpha E_{\theta_0}[T],$$

as in Eqs. (6) and (7).

The inequality  $f_3(t) < k_1 f_1(t) + k_2 f_2(t)$  is equivalent to  $e^{bt} - a_1 - a_2 t < 0$ , where  $a_i = k_i c(\theta_0)/c(\theta_1)$ ,  $i = 1, 2$ , and  $b = \theta_1 - \theta_0$  with a fixed  $\theta_1 \neq \theta_0$ . If the resulting  $\varphi_0(t)$  does not depend on the specific  $\theta_1 \neq \theta_0$ , then it is the desired solution to the problem.

The function  $e^{bt} - a_1 - a_2 t$  has positive second derivative and is therefore convex. Hence  $\{t: e^{bt} - a_1 - a_2 t < 0\}$  is either  $(c, \infty)$ , or  $(-\infty, c')$ , or  $(c_1, c_2)$ . In the first two cases,  $\varphi_0$  is a one-sided test having monotone power, which contradicts the property of  $\beta_{\varphi_0}(\theta) = E_\theta[\varphi_0(T)]$  being minimized at  $\theta_0$ . Hence

$$\varphi_0(t) = \begin{cases} 0 & \text{if } c_1 < t < c_2 \\ \gamma_i & \text{if } t = c_i, \quad i = 1, 2 \\ 1 & \text{if } t < c_1 \text{ or } t > c_2, \end{cases}$$

where  $c_1 < c_2$  and  $0 \leq \gamma_i \leq 1$ ,  $i = 1, 2$  are such that  $E_{\theta_0}[\varphi_0(T)] = \alpha$  and  $E_{\theta_0}[T\varphi_0(T)] = \alpha E_{\theta_0}[T]$ .

To see that  $\varphi_0$  is an unbiased level  $\alpha$  test of  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ , note that  $E_{\theta_0}[\varphi_0(T)] = \alpha$  by construction, and since  $\varphi^*(t) \equiv \alpha$  is in the class in which  $\varphi_0$  maximizes  $\beta(\theta)$  for  $\theta \neq \theta_0$ , we must have  $E_\theta[\varphi_0(T)] \geq E_\theta[\varphi^*(T)] = \alpha$  for all  $\theta \neq \theta_0$ , showing that  $\varphi_0$  is UMP unbiased at level  $\alpha$  for  $H_0$  vs  $H_1$ .

*Remark 6.7.1.* If  $T$  is symmetrically distributed about  $c$  under  $P_{\theta_0}$ , then choosing a test  $\varphi$  so that  $E_{\theta_0}[\varphi(T)] = \alpha$  and  $\varphi$  is symmetric about  $c$ , the condition  $E_{\theta_0}[T\varphi(T)] = \alpha E_{\theta_0}[T]$  is automatically satisfied.

*Proof.* Since  $T$  is symmetrically distributed about  $c$  under  $P_{\theta_0}$ ,

$$\begin{aligned} P_{\theta_0}[T - c < -u] &= P_{\theta_0}[T - c > u] \text{ for all } u \in \mathbb{R}, \text{ so that} \\ E_{\theta_0}[T - c] &= 0, \text{ ie, } E_{\theta_0}[T] = c. \end{aligned}$$

Since  $\varphi$  is symmetric about  $c$ , that is,  $\varphi(c - u) = \varphi(c + u)$  for all  $u \in \mathbb{R}$ , letting  $T^* = T - c$ , we have

$$\begin{aligned} E_{\theta_0}[T\varphi(T)] &= E_{\theta_0}[(T - c)\varphi(T)] + c[\varphi(T)] = E_{\theta_0}[T^*\varphi(T^* + c)] + c\alpha \\ &= \int u\varphi(u + c)p_{T^*}(u; \theta_0) du + \alpha E_{\theta_0}[T] = \alpha E_{\theta_0}[T] \end{aligned}$$

due to symmetry of  $\varphi$  and  $p_{T^*}(\cdot; \theta_0)$ . □

Thus  $\varphi_0$  is obtained by taking  $c_1 = c - k$ ,  $c_2 = c + k$ ,  $\gamma_1 = \gamma_2 = \gamma$ , and choosing  $k$  and  $0 \leq \gamma \leq 1$  so that  $E_{\theta_0}[\varphi_0(T)] = \alpha$ .

### 6.7.2 UMP Unbiased Test for $H_0: \theta_1 \leq \theta \leq \theta_2$ vs $H_1: \theta \notin [\theta_1, \theta_2]$ in a Single-Parameter Exponential Family

Since all tests have continuous power functions in the context of exponential families, we can use [Lemma 6.5.1](#) to narrow down our search among tests which are similar of size  $\alpha$

on the common boundary of  $[\theta_1, \theta_2]$  and  $[\theta_1, \theta_2]^c$  which is  $\omega = \{\theta_1, \theta_2\}$  (ie, among all tests satisfying  $E_{\theta_i}[\varphi(T)] = \alpha, i = 1, 2$ ).

This leads to the problem of maximizing  $\int \varphi(t) f_3(t) dt$  with respect to  $\varphi \in \mathcal{F}$  subject to the conditions  $\int \varphi(t) f_i(t) dt = \alpha, i = 1, 2$ , where  $f_1, f_2$ , and  $f_3$  are given in Eqs. (11) and (12).

Again by the Generalized N-P Lemma, the UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$  in this problem is found to be *of the same form* as the UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ , but here  $c_1, c_2$  and  $0 \leq \gamma_i \leq 1$  are determined by  $E_{\theta_i}[\varphi_0(T)] = \alpha, i = 1, 2$ .

### 6.7.3 Examples

**Example 6.7.1.** Let  $X = (X_1, \dots, X_n)$  be a random sample from  $N(\mu, \sigma^2)$  with  $\sigma$  known. The goal is to find UMP unbiased level  $\alpha$  tests for (a)  $H_0: \mu = \mu_0$  vs  $H_1: \mu \neq \mu_0$ , and (b)  $H_0: \mu_1 \leq \mu \leq \mu_2$  vs  $H_1: \mu \notin [\mu_1, \mu_2]$ .

*Solution.*

(a) Denote  $T = \sqrt{n}(\bar{X} - \mu_0)/\sigma$ . Then the UMP unbiased level  $\alpha$  test is

$$\varphi(\mathbf{x}) = \begin{cases} 0 & \text{if } |T| \leq \Phi^{-1}(1 - \alpha/2) \\ 1 & \text{if } |T| > \Phi^{-1}(1 - \alpha/2), \end{cases}$$

because  $\bar{X}$  is sufficient,  $E_{\mu_0}[\varphi(\mathbf{X})] = \alpha$  and  $E_{\mu_0}[(\bar{X} - \mu_0)\varphi(\mathbf{X})] = 0$  holds since  $\bar{X} - \mu_0$  is symmetrically distributed about 0.

(b) Transform  $(\mu_1, \mu_2)$  to  $(\bar{\mu}, \Delta) = ((\mu_1 + \mu_2)/2, (\mu_2 - \mu_1)/2)$  and let  $\theta = \mu - \bar{\mu}$ . Also let  $Y_i = X_i - \bar{\mu}, i = 1, \dots, n$ . Then the problem can be equivalently expressed as that of testing  $H_0: -\Delta \leq \theta \leq \Delta$  vs  $H_1: |\theta| > \Delta$  based on  $Y_1, \dots, Y_n$  which are iid  $N(\theta, \sigma^2)$ , and  $T = \bar{Y} = n^{-1} \sum_{i=1}^n Y_i$  is sufficient for  $\theta$  in  $(Y_1, \dots, Y_n)$ .

Then the UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$  is

$$\varphi_0(t) = \begin{cases} 0 & \text{if } c_1 < t < c_2 \\ 1 & \text{if } t \leq c_1 \text{ or } t \geq c_2, \end{cases}$$

where  $c_1, c_2$  are such that  $E_{\theta=-\Delta}[\varphi_0(T)] = E_{\theta=\Delta}[\varphi_0(T)] = \alpha$ .

Since  $T \sim N(\theta, \sigma^2/n)$ , choosing  $c_1 = -c$  and  $c_2 = c$ , we have

$$\begin{aligned} E_{\theta=-\Delta}[\varphi_0(T)] &= P_{-\Delta}[T \leq -c] + P_{-\Delta}[T \geq c] \\ &= \Phi(\sqrt{n}(-c + \Delta)/\sigma) + \{1 - \Phi(\sqrt{n}(c + \Delta)/\sigma)\} \\ &= \{1 - \Phi(\sqrt{n}(c - \Delta)/\sigma)\} + \Phi(\sqrt{n}(-c - \Delta)/\sigma) \\ &= P_{\Delta}[T \geq c] + P_{\Delta}[T \leq -c] = E_{\theta=\Delta}[\varphi_0(T)]. \end{aligned}$$

Thus, if we choose  $c_1 = -c$  and  $c_2 = c$  in  $\varphi_0(t)$  and let  $c$  be such that  $E_{\theta=-\Delta}[\varphi_0(T)] = \alpha$ , then  $E_{\theta=\Delta}[\varphi_0(T)] = \alpha$  is automatic.

**Example 6.7.2.** Let  $X = (X_1, \dots, X_n)$  be a random sample from  $\text{Exp}(\theta)$ . We will find the UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .

*Solution.* Here  $T = \sum_{i=1}^n X_i$  is sufficient for  $\theta$  in  $\mathbf{X}$  with  $T \mid \theta_0 \sim \text{Gamma}(n, 1)$  under  $H_0$ . The UMP unbiased level  $\alpha$  test is

$$\varphi(t) = \begin{cases} 0 & \text{if } c_1 < t/\theta_0 < c_2 \\ 1 & \text{if } t/\theta_0 \leq c_1 \text{ or } t/\theta_0 \geq c_2, \end{cases}$$

where  $c_1$  and  $c_2$  are determined by

- (i)  $\int_{c_1}^{c_2} f_n(y) dy = 1 - \alpha$ , and
- (ii)  $\int_{c_1}^{c_2} yf_n(y) dy = n(1 - \alpha)$ , where  $f_n(y) = y^{n-1}e^{-y}/\Gamma(n)$ ,  $y > 0$ , is the pdf of  $\text{Gamma}(n, 1)$  with mean  $n$ . Condition (ii) can be equivalently expressed in any of the following two ways, using integration by part and (i):
- (ii)'  $\int_{c_1}^{c_2} f_{n+1}(y) dy = 1 - \alpha$  since  $yf_n(y) = nf_{n+1}(y)$ , or
- (ii)''  $e^{-c_1}c_1^n = e^{-c_2}c_2^n$  since

$$\int_{c_1}^{c_2} yf_n(y) dy = (e^{-c_1}c_1^n - e^{-c_2}c_2^n)/\Gamma(n) + n(1 - \alpha).$$

For moderately large  $n$  and for  $\theta_0$  neither too large nor too small, CLT provides a reasonable approximation for  $c_1, c_2$  determined by

$$\int_0^{c_1} f_n(y) dy = \int_{c_2}^{\infty} f_n(y) dy = \alpha/2.$$

**Example 6.7.3.** Let  $X$  be a random sample of size 1 from  $\text{Geom}(p)$  with pmf

$$f(x, p) = pq^{x-1}, \quad x = 1, 2, \dots, \text{ where } q = 1 - p.$$

We wish to find UMP unbiased test at level  $\alpha$  for  $H_0: p = p_0$  vs  $H_1: p \neq p_0$ .

*Solution.* The UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$  is

$$\varphi_0(x) = \begin{cases} 0 & \text{if } c_1 < x < c_2 \\ \gamma_i & \text{if } x = c_i, \quad i = 1, 2 \\ 1 & \text{if } x < c_1 \text{ or } x > c_2, \end{cases}$$

where  $c_1, c_2, \gamma_1$ , and  $\gamma_2$  are determined by

$$\begin{aligned} E_{p_0}[1 - \varphi_0(X)] &= p_0 q_0^{-1} \left[ S_{c_2-1} - S_{c_1} + (1 - \gamma_1)q_0^{c_1} + (1 - \gamma_2)q_0^{c_2} \right] \\ &= 1 - \alpha, \text{ and} \\ E_{p_0}[X\{1 - \varphi_0(X)\}] &= p_0 q_0^{-1} \left[ S_{c_2-1}^* - S_{c_1}^* + (1 - \gamma_1)c_1 q_0^{c_1} + (1 - \gamma_2)c_2 q_0^{c_2} \right] \\ &= (1 - \alpha)p_0^{-1}, \end{aligned}$$

with

$$S_r = \sum_{i=1}^r q_0^i = q_0 p_0^{-1} (1 - q_0^r) \text{ and } S_r^* = \sum_{i=1}^r i q_0^i = q_0 p_0^{-2} [1 - (1 + r p_0) q_0^r].$$

## 6.8 Locally Best Tests

We have seen so far that UMP tests exist for one-sided problems of testing  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  (or vice versa) in MLR families, and UMP unbiased tests exist for two-sided problems of testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  or for testing  $H_0: \theta_1 \leq \theta \leq \theta_2$  vs  $H_1: \theta \notin [\theta_1, \theta_2]$  in single-parameter exponential families. In subsequent developments, we are also going to construct UMP unbiased tests for one-sided and two-sided problems concerning any one parameter in multiparameter exponential families when other parameters (called nuisance parameters) are unknown but are of no concern. However, UMP or UMP unbiased tests are not available outside the MLR or the exponential families, so in more general situations, we would have to lower our expectation and settle for some more modest criterion of optimality.

In the one-sided problem, suppose that for every test  $\varphi$ , the power function  $\beta_\varphi(\theta)$  has a continuous derivative which can be obtained by differentiating under the integral. In the context of a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $f(x; \theta)$ , this means

$$\begin{aligned}\beta'_\varphi(\theta) &= \frac{d}{d\theta} E_\theta[\varphi(\mathbf{X})] = \frac{d}{d\theta} \int \varphi(\mathbf{x}) f(\mathbf{x}; \theta) d\mathbf{x} \\ &= \int \varphi(\mathbf{x}) \frac{\partial f(\mathbf{x}; \theta)}{\partial \theta} d\mathbf{x} \\ &= \int \varphi(\mathbf{x}) \left( \sum_{i=1}^n \frac{\partial \log f(\mathbf{x}; \theta)}{\partial \theta} \right) f(\mathbf{x}; \theta) d\mathbf{x} \\ &= E_\theta \left[ \varphi(\mathbf{X}) \sum_{i=1}^n \dot{l}(X_i; \theta) \right],\end{aligned}$$

where  $\dot{l}(x; \theta) = \frac{\partial \log f(x; \theta)}{\partial \theta}$  is the partial derivative of the log likelihood  $l(x; \theta) = \log f(x; \theta)$  with respect to  $\theta$  as in [Section 5.2.1](#).

We shall use the other notations introduced in [Section 5.2.1](#), assuming that the regularity conditions introduced there holds here and use the results obtained under those conditions.

We now formulate our criterion for local optimality in the one-sided problem in a single-parameter family.

**Definition 6.8.1.** A test  $\varphi_0$  is said to be a *locally most powerful* (LMP) test at level  $\alpha$  for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  if

- (i)  $\beta_{\varphi_0}(\theta_0) = E_{\theta_0}[\varphi(\mathbf{X})] = \alpha$ , and
- (ii)  $\beta'_{\varphi_0}(\theta_0) \geq \beta'_{\varphi}(\theta_0)$ , that is,

$$E_{\theta_0} \left[ \varphi_0(\mathbf{X}) \sum_{i=1}^n \dot{l}(X_i; \theta_0) \right] \geq E_{\theta_0} \left[ \varphi(\mathbf{X}) \sum_{i=1}^n \dot{l}(X_i; \theta_0) \right]$$

for all tests  $\varphi$  satisfying (i).

Applying the Generalized N-P Lemma with  $m = 1$ ,  $f_1(\mathbf{x}) = f(\mathbf{x}; \theta_0)$  and  $f_2(\mathbf{x}) = \left\{ \sum_{i=1}^n \dot{l}(x_i; \theta_0) \right\} f(\mathbf{x}; \theta_0)$ , we see that any test of the form:

$$\varphi_0(\mathbf{x}) = \begin{cases} 0 & \text{if } f_2(\mathbf{x}) < kf_1(\mathbf{x}) \text{ ie, } \sum_{i=1}^n \dot{l}(x_i; \theta_0) < k \\ \gamma & \text{if } f_2(\mathbf{x}) = kf_1(\mathbf{x}) \\ 1 & \text{if } f_2(\mathbf{x}) > kf_1(\mathbf{x}) \end{cases}$$

is an LMP test at level  $\alpha$  for  $H_0$  vs  $H_1$  provided that

$$E_{\theta_0}[\varphi_0(\mathbf{X})] = P_{\theta_0} \left[ \sum_{i=1}^n \dot{l}(X_i; \theta_0) > k \right] + \gamma P_{\theta_0} \left[ \sum_{i=1}^n \dot{l}(X_i; \theta_0) = k \right] = \alpha.$$

For large  $n$ , we can find the approximate value of  $k$  for a given  $\alpha$  by the CLT. Recall that under the regularity conditions in [Section 5.2.1](#),

$$E_{\theta_0}[\dot{l}(X_i; \theta_0)] = 0 \quad \text{and} \quad \text{Var}_{\theta_0}[\dot{l}(X_i; \theta_0)] = I(\theta_0),$$

where  $I(\theta_0)$  is the Fisher-information. Thus

$$n^{-1/2} \sum_{i=1}^n \dot{l}(X_i; \theta_0) \xrightarrow{\mathcal{L}} N(0, I(\theta_0)) \text{ under } P_{\theta_0}.$$

Therefore, for large  $n$ , the critical value  $k$  for a given  $\alpha$  can be approximated by

$$k = k_{n,\alpha} \simeq \sqrt{nI(\theta_0)} \Phi^{-1}(1 - \alpha).$$

For testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ , we assume that the power function  $\beta_\varphi(\theta)$  of every test  $\varphi$  has two continuous derivatives which can be obtained by differentiating under the integral. Thus for a test  $\varphi$  based on a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $f(x; \theta)$ , we have  $\beta'_\varphi(\theta) = E_\theta[\varphi(\mathbf{X}) \sum_{i=1}^n \ddot{l}(X_i; \theta)]$  as before, and

$$\begin{aligned} \beta''_\varphi(\theta) &= \int \varphi(\mathbf{x}) \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta^2} d\mathbf{x} \\ &= E_\theta \left[ \varphi(\mathbf{X}) \left\{ \sum_{i=1}^n \ddot{l}(X_i; \theta) + \left( \sum_{i=1}^n \dot{l}(X_i; \theta) \right)^2 \right\} \right], \end{aligned}$$

using the identity:

$$\begin{aligned} \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta^2} &= \left[ \frac{\partial^2 \log f(\mathbf{x}; \theta)}{\partial \theta^2} + \left( \frac{\partial \log f(\mathbf{x}; \theta)}{\partial \theta} \right)^2 \right] \\ &= \left[ \sum_{i=1}^n \ddot{l}(x_i; \theta) + \left( \sum_{i=1}^n \dot{l}(x_i; \theta) \right)^2 \right] f(\mathbf{x}; \theta), \end{aligned}$$

where  $\dot{l}(x; \theta)$  and  $\ddot{l}(x; \theta)$  are the first two partial derivatives of  $l(x; \theta) = \log f(x; \theta)$  with respect to  $\theta$ .

**Definition 6.8.2.** A test  $\varphi_0$  is said to be an LMP (locally) unbiased test at level  $\alpha$  for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  if

- (i)  $\beta_{\varphi_0}(\theta_0) = E_{\theta_0}[\varphi_0(\mathbf{X})] = \alpha$ ,
- (ii)  $\beta'_{\varphi_0}(\theta_0) = E_{\theta_0}\left[\varphi_0(\mathbf{X}) \sum_{i=1}^n \dot{l}(X_i; \theta_0)\right] = 0$ , and
- (iii)  $\beta''_{\varphi_0}(\theta_0) \geq \beta''_\varphi(\theta_0)$ , that is,

$$\begin{aligned} & E_{\theta_0}\left[\varphi_0(\mathbf{X}) \left\{ \sum_{i=1}^n \ddot{l}(X_i; \theta_0) + \left( \sum_{i=1}^n \dot{l}(X_i; \theta_0) \right)^2 \right\} \right] \\ & \geq E_{\theta_0}\left[\varphi(\mathbf{X}) \left\{ \sum_{i=1}^n \ddot{l}(X_i; \theta_0) + \left( \sum_{i=1}^n \dot{l}(X_i; \theta_0) \right)^2 \right\} \right] \end{aligned}$$

for all tests  $\varphi$  satisfying (i) and (ii).

We can now apply the Generalized N-P Lemma with  $m = 2$ ,  $f_1(\mathbf{x}) = f(\mathbf{x}; \theta_0)$ ,  $f_2(\mathbf{x}) = Uf(\mathbf{x}; \theta_0)$ ,  $f_3(\mathbf{x}) = Vf(\mathbf{x}; \theta_0)$ , where

$$U = \sum_{i=1}^n \dot{l}(x_i; \theta_0) \quad \text{and} \quad V = \sum_{i=1}^n \ddot{l}(x_i; \theta_0) + \left( \sum_{i=1}^n \dot{l}(x_i; \theta_0) \right)^2,$$

to see that any test of the form:

$$\varphi_{k_1, k_2}(\mathbf{x}) = \begin{cases} 0 & \text{if } V < k_1 + k_2 U \\ \gamma & \text{if } V = k_1 + k_2 U \\ 1 & \text{if } V > k_1 + k_2 U \end{cases}$$

is an LMP unbiased test for  $H_0$  vs  $H_1$  at level  $\alpha$ , provided that  $k_1$  and  $k_2$  are chosen so as to satisfy (i) and (ii).

Finding  $k_1$ ,  $k_2$  for a given  $\alpha$  is difficult in general, but in the special case when  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$  under  $P_{\theta_0}$ , the problem is simplified because we can take  $k_2 = 0$ .

**Proposition 6.8.1.** If  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$  under  $P_{\theta_0}$ , then for any  $k_1$ , condition (ii) holds iff  $k_2 = 0$ .

*Proof.* If  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$  under  $P_{\theta_0}$ , then for  $k_2 = 0$ ,

$$\begin{aligned} E_{\theta_0}[U\varphi_{k_1, 0}(\mathbf{X})] &= E_{\theta_0}[UI_{(k_1, \infty)}(V)] = E_{\theta_0}[-UI_{(k_1, \infty)}(V)] \\ &= -E_{\theta_0}[U\varphi_{k_1, 0}(\mathbf{X})]. \end{aligned}$$

Hence  $E_{\theta_0}[U\varphi_{k_1, 0}(\mathbf{X})] = 0$ . Conversely, if  $k_2 > 0$ , then

$$\begin{aligned} E_{\theta_0}[U\varphi_{k_1, k_2}(\mathbf{X})] &= \int_{u=-\infty}^{\infty} \int_{v=k_1+k_2 u}^{\infty} uf(u, v; \theta_0) du dv \\ &= \sum_{j=1}^3 \iint_{S_j} uf(u, v; \theta_0) du dv \\ &= I_1 + I_2 + I_3, \text{ say} \end{aligned}$$

where

$$\begin{aligned} S_1 &= \{(u, v): u \geq 0, k_1 + k_2 u \leq v < \infty\}, \\ S_2 &= \{(u, v): u < 0, k_1 - k_2 u \leq v < \infty\}, \text{ and} \\ S_3 &= \{(u, v): u < 0, k_1 + k_2 u \leq v < k_1 - k_2 u\}. \end{aligned}$$

Since  $(u, v) \in S_1 \iff (-u, v) \in S_2$  and  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$  under  $\theta_0$ , it follows that

$$\begin{aligned} I_2 &= \iint_{S_2} uf(u, v; \theta_0) du dv = - \iint_{S_2} (-u)f(-u, v; \theta_0) du dv \\ &= - \iint_{S_1} uf(u, v; \theta_0) du dv = -I_1. \end{aligned}$$

Hence

$$E_{\theta_0}[U\varphi_{k_1, k_2}(\mathbf{X})] = I_3 = \iint_{S_3} uf(u, v; \theta_0) du dv < 0$$

because  $u < 0$  on  $S_3$ . Hence in order to satisfy  $E[U\varphi_{k_1, k_2}(\mathbf{X})] = 0$ , we must have  $P_{\theta_0}[S_3] = 0$ . However, if  $P_{\theta_0}[S_3] = 0$ , then  $E_{\theta_0}[\varphi_{k_1, k_2}(\mathbf{X})] = E_{\theta_0}[\varphi_{k_1, 0}(\mathbf{X})]$ . Thus both (i) and (ii) are satisfied with  $k_2 = 0$ . The case of  $k_2 < 0$  is treated in the same way.  $\square$

**Example 6.8.1.** Let  $X = (X_1, \dots, X_n)$  be a random sample from a logistic distribution with location parameter  $\theta$ , that is, from

$$f(x, \theta) = \exp[-(x - \theta)] / \{1 + \exp[-(x - \theta)]\}^2, \quad -\infty < x < \infty.$$

We wish to find the locally best

- (a) test for  $H_0: \theta = 0$  vs  $H_1: \theta > 0$  at level  $\alpha$ , and
- (b) unbiased test for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  at level  $\alpha$ .

*Solution.* Note that

$$\begin{aligned} l(x, \theta) &= \log f(x, \theta) = -(x - \theta) - 2 \log[1 + \exp\{-(x - \theta)\}], \\ \dot{l}(x, \theta) &= \frac{\partial}{\partial \theta} l(x, \theta) = [\exp(x - \theta) - 1]/[\exp(x - \theta) + 1], \\ \ddot{l}(x, \theta) &= \frac{\partial^2}{\partial \theta^2} l(x, \theta) = -2 \exp(x - \theta)/(\exp(x - \theta) + 1)^2. \end{aligned}$$

Thus we have

$$\begin{aligned} U &= \sum_{i=1}^n \dot{l}(X_i, 0) = \sum_{i=1}^n (\exp^{X_i} - 1)/(\exp^{X_i} + 1), \\ V &= \sum_{i=1}^n \ddot{l}(X_i, 0) + \left\{ \sum_{i=1}^n \dot{l}(X_i, 0) \right\}^2 = -2 \sum_{i=1}^n \exp^{X_i}/(\exp^{X_i} + 1)^2 + U^2. \end{aligned}$$

(a) The LMP test at level  $\alpha$  for  $H_0: \theta = 0$  vs  $H_1: \theta > 0$  is

$$\varphi_1(\mathbf{x}) = \begin{cases} 0 & \text{if } U < k \\ \gamma & \text{if } U = k \\ 1 & \text{if } U > k, \end{cases}$$

where  $P_0[U > k] + \gamma P_0[U = k] = \alpha$ .

(b) The LMP (locally) unbiased test at level  $\alpha$  for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  is

$$\varphi_0(\mathbf{x}) = \begin{cases} 0 & \text{if } V < k_1 + k_2 U \\ \gamma & \text{if } V = k_1 + k_2 U \\ 1 & \text{if } V > k_1 + k_2 U, \end{cases}$$

where  $0 \leq \gamma \leq 1$  and  $k_1, k_2$  are determined by (i)  $E_0[\varphi_0(\mathbf{X})] = \alpha$  and (ii)  $E_0[\varphi_0(\mathbf{X})U] = 0$ .

However, by virtue of symmetry of  $f(x, 0)$ ,  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$ . Hence  $k_2$  can be taken to be equal to 0 and it is enough to choose  $0 \leq \gamma \leq 1$  and  $k_1$  to satisfy (i), since (ii) automatically holds.

The condition  $(U, V) \stackrel{\mathcal{D}}{=} (-U, V)$  under  $P_{\theta_0}$  is satisfied in a number of situations such as

- (i)  $X \sim \text{Cauchy}(\theta, 1)$  with  $f(x; \theta) = \frac{1}{\pi} \frac{1}{1+(x-\theta)^2}$ ,  $x \in \mathbb{R}$ ,
- (ii)  $X \sim \text{Double Exponential}(\theta)$  with  $f(x; \theta) = \frac{1}{2} \exp[-|x - \theta|]$ ,  $x \in \mathbb{R}$ , etc.

## 6.9 UMP Unbiased Tests in the Presence of Nuisance Parameters: Similarity and Completeness

In many situations, the distribution of the observed  $X$  depends on several parameters and we want to test a null hypothesis  $H_0$  against an alternative  $H_1$  concerning only one of these parameters. The other parameters are called *nuisance parameters*.

**Example 6.9.1.** Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $N(\mu, \sigma^2)$ , where both  $\mu, \sigma^2$  are unknown. We want to test  $H_1: \mu = 0$  vs  $H_1: \mu \neq 0$ . Here  $\sigma^2$  is a nuisance parameter.

**Example 6.9.2.** Let  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  be independent random samples from  $\text{Poisson}(\lambda)$  and  $\text{Poisson}(\mu)$ , respectively. We want to test  $H_0: \mu \leq a\lambda$  vs  $H_1: \mu > a\lambda$  for a given  $a$ . To see how this problem involves testing for one parameter in the presence of a nuisance parameter, first consider the joint distribution of the sufficient statistic  $(X, Y) = (\sum_{i=1}^m X_i, \sum_{i=1}^n Y_i)$  given by:

$$f_{X,Y}(x, y; \lambda, \mu) = e^{-(m\lambda+n\mu)} \frac{(m\lambda)^x (n\mu)^y}{x! y!}, \quad x = 0, 1, 2, \dots, y = 0, 1, 2, \dots$$

Now reparametrize:  $\pi = n\mu/(m\lambda + n\mu)$ ,  $\xi = m\lambda + n\mu$  and transform the data:  $U = Y$ ,  $T = X + Y$ . Then the joint distribution of  $(U, T)$  is

$$f_{U,T}(u, t; \pi, \xi) = \left\{ \binom{t}{u} \pi^u (1-\pi)^{t-u} \right\} \left\{ e^{-\xi} \xi^t / t! \right\},$$

$t = 0, 1, 2, \dots$  and  $u = 0, 1, \dots, t$ . Now  $H_0: \mu \leq a\lambda \iff \pi \leq na/(m + na)$  and  $H_1: \mu > a\lambda \iff \pi > na/(m + na)$  while  $\xi$  is a nuisance parameter.

Let  $\{P_\theta^X, \theta \in \Theta\}$  be a family of probabilities on  $(\mathfrak{X}, \mathcal{A})$ , where  $\Theta \subset \mathbb{R}^{k+1}$  and  $\theta$  is a  $(k+1)$ -dim vector denoted by

$$\begin{aligned}\boldsymbol{\theta} &= (\theta_1, \theta_2, \dots, \theta_{k+1}) := (\theta, \tau_1, \dots, \tau_k) := (\theta, \boldsymbol{\tau}) \text{ with} \\ \theta &= \theta_1 \text{ and } \boldsymbol{\tau} = (\theta_2, \dots, \theta_{k+1}).\end{aligned}$$

Based on a random element  $X$  generated by an unknown element  $P_\theta^X$  of this family, we want to find unbiased level  $\alpha$  tests in three situations where the null and the alternative hypotheses are given as:

**Problem 1.**  $H_{10}: \theta \leq \theta_0$  vs  $H_{11}: \theta > \theta_0$  (testing  $\theta \geq \theta_0$  vs  $\theta < \theta_0$  is analogous),

**Problem 2.**  $H_{20}: \theta = \theta_0$  vs  $H_{21}: \theta \neq \theta_0$ ,

**Problem 3.**  $H_{30}: \theta_1 \leq \theta \leq \theta_2$  vs  $H_{31}: \theta \notin [\theta_1, \theta_2]$ ,

in each case,  $\boldsymbol{\tau}$  being the nuisance parameter.

A more careful description of the null and the alternative hypotheses in *these three problems* are

$$\begin{aligned}H_{i0} &: (\theta, \boldsymbol{\tau}) \in \Theta_{i0} \text{ vs } H_{i1} : (\theta, \boldsymbol{\tau}) \in \Theta_{i1}, \text{ where} \\ \Theta_{10} &= \{(\theta, \boldsymbol{\tau}): \theta \leq \theta_0, \boldsymbol{\tau} \in \Omega\}, \quad \Theta_{11} = \{(\theta, \boldsymbol{\tau}): \theta > \theta_0, \boldsymbol{\tau} \in \Omega\}, \\ \Theta_{20} &= \{(\theta, \boldsymbol{\tau}): \theta = \theta_0, \boldsymbol{\tau} \in \Omega\}, \quad \Theta_{21} = \{(\theta, \boldsymbol{\tau}): \theta \neq \theta_0, \boldsymbol{\tau} \in \Omega\}, \\ \Theta_{30} &= \{(\theta, \boldsymbol{\tau}): \theta_1 \leq \theta \leq \theta_2, \boldsymbol{\tau} \in \Omega\}, \\ \Theta_{31} &= \{(\theta, \boldsymbol{\tau}): \theta \notin [\theta_1, \theta_2], \boldsymbol{\tau} \in \Omega\}, \text{ where } \Omega \subset \mathbb{R}^k.\end{aligned}$$

Now in each of these cases, the null hypothesis and the alternative hypothesis sets  $\Theta_{i0}$  and  $\Theta_{i1}$  have a common boundary  $\omega_i$ , where

$$\omega_1 = \omega_2 = \{\theta_0\} \times \Omega \quad \text{and} \quad \omega_3 = \{\theta_1, \theta_2\} \times \Omega.$$

In the sequel, we assume that the power functions  $\beta_\varphi(\theta, \boldsymbol{\tau})$  of all tests  $\varphi$  are continuous in  $(\theta, \boldsymbol{\tau})$ , so that all unbiased level  $\alpha$  tests are similar tests of size  $\alpha$  on the boundary  $\omega_i$  of  $\Theta_{i0}$  and  $\Theta_{i1}$ ,  $i = 1, 2, 3$  in each of the three problems under consideration.

We now omit the subscript  $i$  from the boundary  $\omega_i$  of  $\Theta_{i0}$  and  $\Theta_{i1}$ , and continue the development in some generality. Restrict attention to the family of distributions  $\mathcal{P}_\omega^X = \{P_{\theta, \boldsymbol{\tau}}^X, (\theta, \boldsymbol{\tau}) \in \omega\}$  and suppose that  $T$  is a sufficient statistic for the family  $\mathcal{P}_\omega^X$  (ie, the conditional distribution of  $X$  given  $T$  is the same for all  $(\theta, \boldsymbol{\tau}) \in \omega$ ). Indeed, we should be looking for a sufficient statistic  $T$  for  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_k)$  since  $\theta_0$  in Problems 1, 2 and  $\theta_1, \theta_2$  in Problem 3 are fixed and known.

**Definition 6.9.1** (Tests of Neyman Structure). Tests satisfying  $E[\varphi(X)|T = t] = \alpha$  for almost all  $t$  under  $\mathcal{P}_\omega^T$  (ie, for all  $t \notin N$ , where  $P_\theta[T \in N] = 0$  for all  $\theta \in \omega$ ) are called tests of Neyman-structure with respect to the sufficient statistic  $T$ .

If a test is of Neyman-structure, then

$$E_\theta[\varphi(X)] = E_\theta E[\varphi(X)|T] = E_\theta[\alpha] = \alpha \quad \text{for all } \theta \in \omega,$$

that is,  $\varphi$  is similar of size  $\alpha$  on  $\omega$ . Thus all tests of Neyman-structure are similar on  $\omega$ .

A sufficient statistic  $T$  for  $\mathcal{P}_\omega^X$  is *complete* if  $\int g(t)f_T(t; \theta) dt = 0$  for all  $\theta \in \omega \Rightarrow g(t) = 0$  for almost all  $t$  under  $\mathcal{P}_\omega^T$  (ie,  $g(t) = 0$  for all  $t \notin N$  where  $P_\theta[T \in N] = 0$  for all  $\theta \in \omega$ ). The property of completeness of a sufficient statistic has already been used in connection with UMVUE in [Section 5.1](#).

A sufficient statistic  $T$  for  $\mathcal{P}_\omega^X$  is *boundedly complete* if  $\int g(t)f_T(t; \theta) dt = 0$  for all  $\theta \in \omega$ ,  $g$  bounded  $\Rightarrow g(t) = 0$  for almost all  $t$  under  $\mathcal{P}_\omega^T$ . [If  $T$  is complete, then it is boundedly complete, but the converse is not true.]

We now have the following theorem.

**Theorem 6.9.1.** *Suppose that there exists a sufficient statistic  $T$  for  $\mathcal{P}_\omega^X$ . If  $T$  is boundedly complete, then every test which is similar on  $\omega$ , is of Neyman-structure with respect to  $T$ .*

*Proof.* Let  $\varphi$  be a similar test of size  $\alpha$  on  $\omega$ . Then

$$\begin{aligned} E_\theta[\varphi(X)] = \alpha &\quad \text{for all } \theta \in \omega \Leftrightarrow \\ E_\theta[\varphi(X) - \alpha] = E_\theta E[\varphi(X) - \alpha|T] &= 0 \quad \text{for all } \theta \in \omega. \end{aligned}$$

Since  $\psi(t) = E[\varphi(X) - \alpha|T = t]$  is a bounded function of  $t$ ,

$$E_\theta[\psi(t) = 0] \quad \text{for all } \theta \in \omega$$

$\Rightarrow \psi(t) = 0$  for almost all  $t$  under  $\mathcal{P}_\omega^T \Leftrightarrow E[\varphi(X)|T = t] = \alpha$  for almost all  $t$  under  $\mathcal{P}_\omega^T$ , that is,  $\varphi$  is a test of Neyman-structure with respect to  $T$ .  $\square$

The converse (ie, if  $T$  is not boundedly complete, then there exists a test which is similar on  $\alpha$  but is not of Neyman-structure with respect to  $T$ ) is also true (see [\[3, p. 134\]](#)), but we do not need it for our purpose.

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from a  $(k + 1)$ -parameter exponential family with pdf/pmf:

$$f_X(x; \theta, \tau) = C_0(\theta, \tau) \exp \left[ \theta U(x) + \sum_{j=1}^k \tau_j T_j(x) \right] h_0(x), \quad (\theta, \tau) \in \Theta,$$

where  $\Theta$  is a convex set in  $\mathbb{R}^{k+1}$ .

We shall now use the above theorem to find UMP unbiased level  $\alpha$  tests in the three problems listed above.

Restrict attention to tests based on

$$U = \sum_{i=1}^n U(X_i) \quad \text{and} \quad \mathbf{T} = (T_1, \dots, T_k) = \left( \sum_{i=1}^n T_1(X_i), \dots, \sum_{i=1}^n T_k(X_i) \right),$$

since  $(U, \mathbf{T})$  is sufficient for  $(\theta, \tau)$  in  $X$ . The joint pdf/pmf of  $(U, \mathbf{T})$  is of the form:

$$f_{U,\mathbf{T}}(u, \mathbf{t}; \theta, \tau) = C(\theta, \tau) \exp \left[ \theta u + \sum_{j=1}^k \tau_j t_j \right] h(u, \mathbf{t}), \quad (\theta, \tau) \in \Theta.$$

We assume that:

- (a)  $\Theta = I \times \Omega$  where  $I$  is an interval in  $\mathbb{R}$  and  $\Omega$  is a convex set in  $\mathbb{R}^k$ , and recall that the boundaries of the null and the alternative hypotheses in the three problems are

$$\omega_1 = \omega_2 = \{(\theta_0, \tau): \tau \in \Omega\} \text{ in Problems 1 and 2, and}$$

$$\omega_3 = \omega_{31} \cup \omega_{32} = \{(\theta_1, \tau): \tau \in \Omega\} \cup \{(\theta_2, \tau): \tau \in \Omega\} \text{ in Problem 3.}$$

- (b) In Problems 1 and 2,  $\theta_0$  is an interior point of  $I$  and in Problem 3,  $\theta_1$  and  $\theta_2$  are interior points of  $I$ .  
(c)  $\Omega$  contains a nondegenerate  $k$ -dim rectangle in  $\mathbb{R}^k$ .

In each of the three problems under consideration, the statistic  $\mathbf{T}$  is sufficient for  $\mathcal{P}_\omega^X$  on the common boundary  $\omega$  of the null and the alternative hypotheses (ie, for  $\mathcal{P}_{\omega_i}^X$  in Problems 1, 2, and 3). Also, the marginal distribution of  $\mathbf{T}$  and the conditional distribution of  $U$  given  $\mathbf{T} = \mathbf{t}$  are described by their pdf's/pmf's which are of the form:

$$f_{\mathbf{T}}(\mathbf{t}; \theta, \tau) = C(\theta, \tau) \exp \left[ \sum_{j=1}^k \tau_j t_j \right] h_\theta(\mathbf{t}), \text{ and}$$

$$f_{U|\mathbf{T}}(u; \theta) = C_{\mathbf{t}}(\theta) \exp[\theta u] h_{\mathbf{t}}(u).$$

[In the discrete case,  $h_\theta(\mathbf{t}) = \sum_u \exp[\theta u] h(u, \mathbf{t})$ ,  $C_{\mathbf{t}}(\theta) = 1/h_\theta(\mathbf{t})$ ,  $h_{\mathbf{t}}(u) = h(u, \mathbf{t})$ .]

Note that in all three cases, the sufficient statistic  $\mathbf{T}$  for  $\mathcal{P}_\omega^X$  is complete by virtue of Assumption (c) and is, therefore, boundedly complete. Hence all tests which are similar on  $\omega_i$ ,  $i = 1, 2, 3$  are of Neyman-structure with respect to  $\mathbf{T}$ , which suggests the following strategy:

- (i) For each  $\mathbf{t}$ , restrict attention to the family  $\{f_{U|\mathbf{T}}(u; \theta), \theta \in I\}$  and find the UMP level  $\alpha$  test for  $H_{10}: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  in Problem 1, and the UMP unbiased level  $\alpha$  tests for  $H_{20}: \theta = \theta_0$  vs  $H_{21}: \theta \neq \theta_0$  and for  $H_{30}: \theta_1 \leq \theta \leq \theta_2$  vs  $H_{31}: \theta \notin [\theta_1, \theta_2]$ , respectively, in Problems 2 and 3. [We have already obtained the solutions to these problems in Sections 6.4, 6.7.1, and 6.7.2.]
- (ii) Put these piecewise optimal tests together and show that the resulting test is an unbiased level  $\alpha$  test in each problem.

We now go through the details of each problem.

**Problem 1** (UMP Unbiased Level  $\alpha$  Test for  $H_{10}: \theta \leq \theta_0$  vs  $H_{11}: \theta > \theta_0$ ). We want a test  $\varphi$  such that

- (i)  $E_{\theta, \tau}[\varphi(U, \mathbf{T}) | \mathbf{T} = \mathbf{t}] = \alpha$  for all  $(\theta, \tau) \in \omega_1$ , ie,  $E_{\theta_0}[\varphi(U, \mathbf{T}) | \mathbf{T} = \mathbf{t}] = \alpha$  for almost all  $\mathbf{t}$  under  $\mathcal{P}_{\omega_1}^T$  (written as: a.e.  $\mathcal{P}_{\omega_1}^T$ ), and

(ii)  $E_{\theta, \tau}[\varphi(U, \mathbf{T})] \geq E_{\theta, \tau}[\psi(U, \mathbf{T})]$  for all  $\theta > \theta_0$  and  $\tau \in \Omega$  and for every test  $\psi$  satisfying (i).

[(i) is the condition of Neyman-structure with respect to  $\mathbf{T}$  and (ii) ensures the UMP property in this class.]

For every  $\mathbf{t}$ ,  $\{f_{U|\mathbf{t}}(u; \theta) = C_{\mathbf{t}}(\theta) \exp[\theta u] h_{\mathbf{t}}(u), \theta \in I\}$  is an MLR family, so invoking the result in [Section 6.4](#), the piecewise solution to the problem is seen to be

$$\varphi_1(u, \mathbf{t}) = \begin{cases} 0 & \text{if } u < c_1(\mathbf{t}) \\ \gamma_1(\mathbf{t}) & \text{if } u = c_1(\mathbf{t}) \\ 1 & \text{if } u > c_1(\mathbf{t}), \end{cases}$$

where  $c_1(\mathbf{t})$  and  $0 \leq \gamma_1(\mathbf{t}) \leq 1$  are determined by  $E_{\theta_0}[\varphi_1(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha$ .

To show that  $\varphi_1$  is a UMP level  $\alpha$  test for  $H_{10}: \theta \leq \theta_0$  vs  $H_{11}: \theta > \theta_0$ , it is enough to check that  $\varphi_1$  is a level  $\alpha$  test for  $H_{10}$  vs  $H_{11}$ , using [Lemma 6.5.1](#). By the monotone power property of  $\varphi_1$  in the conditional problem (using part (ii) of [Theorem 6.4.1](#)) and the MLR property of  $\{f_{U|\mathbf{t}}(u; \theta), \theta \in I\}$ ,

$$E_{\theta_1}[\varphi_1(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] \leq E_{\theta_0}[\varphi_1(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha, \text{ a.e. } \mathcal{P}_{\omega_1}^{\mathbf{T}}$$

Hence  $E_{\theta_1, \tau}[\varphi_1(U, \mathbf{T})] \leq \alpha$  for all  $\theta_1 \leq \theta_0$ , showing that  $\varphi_1$  is a level  $\alpha$  test for  $H_{10}$  vs  $H_{11}$ .

**Problem 2** (UMP Unbiased Level  $\alpha$  Test for  $H_{20}: \theta = \theta_0$  vs  $H_{21}: \theta \neq \theta_0$ ). We want a test  $\varphi$  such that

- (i)  $E_{\theta, \tau}[\varphi(U, \mathbf{T}) | \mathbf{T} = \mathbf{t}] = \alpha$  for all  $(\theta, \tau) \in \omega_2$ , ie,  $E_{\theta_0}[\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha$ , a.e.  $\mathcal{P}_{\omega_2}^{\mathbf{T}}$  (as in [Problem 1](#)),
- (ii)  $\frac{\partial}{\partial \theta} E_{\theta, \tau}[\varphi(U, \mathbf{T})] = 0$  on  $\omega_2$ , ie,  $\left. \frac{\partial}{\partial \theta} E_{\theta, \tau}[\varphi(U, \mathbf{T})] \right|_{\theta=\theta_0, \tau \in \Omega} = 0$ , and
- (iii)  $E_{\theta, \tau}[\varphi(U, \mathbf{T})] \geq E_{\theta, \tau}[\psi(U, \mathbf{T})]$  for all  $\theta \neq \theta_0$  and  $\tau \in \Omega$ , whenever  $\psi$  is a test satisfying (i) and (ii).

[(i) is the condition of Neyman-structure with respect to  $\mathbf{T}$ , (ii) is implied by unbiasedness, and (iii) ensures the UMP property in this class.]

Since  $\theta_0$  is an interior point of  $I$ ,  $\frac{\partial}{\partial \theta} E_{\theta, \tau}[\varphi(U, \mathbf{T})]$  exists on  $\{\theta_0\} \times \Omega$  and can be calculated by differentiating under the integral. Thus

$$\begin{aligned} & \frac{\partial}{\partial \theta} E_{\theta, \tau}[\varphi(U, \mathbf{T})] \\ &= \frac{\partial}{\partial \theta} \int \varphi(u, \mathbf{t}) C(\theta, \tau) e^{\theta u + \sum_{j=1}^k \tau_j t_j} h(u, \mathbf{t}) du d\mathbf{t} \\ &= \int \varphi(u, \mathbf{t}) \left[ \frac{\partial C(\theta, \tau)}{\partial \theta} e^{\theta u + \sum_{j=1}^k \tau_j t_j} + C(\theta, \tau) u e^{\theta u + \sum_{j=1}^k \tau_j t_j} \right] h(u, \mathbf{t}) du d\mathbf{t} \\ &= \frac{\frac{\partial C(\theta, \tau)}{\partial \theta}}{C(\theta, \tau)} E_{\theta, \tau}[\varphi(U, \mathbf{T})] + E_{\theta, \tau}[U \varphi(U, \mathbf{T})], \end{aligned}$$

so condition (ii) amounts to the last expression evaluated at  $\theta = \theta_0$  being 0 for all  $\tau \in \Omega$ .

Since  $\varphi^*(u, \mathbf{t}) \equiv \alpha$  is an unbiased level  $\alpha$  test for  $H_{20}$  vs  $H_{21}$ , it follows that

$$\frac{\frac{\partial C(\theta_0, \tau)}{\partial \theta}}{C(\theta_0, \tau)} \alpha + \alpha E_{\theta_0, \tau}[U] = 0, \text{ ie, } \frac{\frac{\partial C(\theta_0, \tau)}{\partial \theta}}{C(\theta_0, \tau)} = -E_{\theta_0, \tau}[U].$$

Condition (ii) thus becomes  $E_{\theta_0, \tau}[U\varphi(U, \mathbf{T}) - \alpha U] = 0$  for all  $\tau \in \Omega$ , that is,  $\int E_{\theta_0}[U\varphi(U, \mathbf{t}) - \alpha U | \mathbf{T} = \mathbf{t}] f_{\mathbf{T}}(\mathbf{t}; \theta_0, \tau) d\mathbf{t} = 0$ , which implies:

$$E_{\theta_0}[U\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha E_{\theta_0}[U | \mathbf{T} = \mathbf{t}], \text{ a.e. } \mathcal{P}_{\omega_2}^T,$$

since  $\mathbf{T}$  is complete.

Conditions (i) and (ii) thus become

$$E_{\theta_0}[\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha \text{ and } E_{\theta_0}[U\varphi(U, \mathbf{t})] = \alpha E_{\theta_0}[U | \mathbf{T} = \mathbf{t}], \text{ a.e. } \mathcal{P}_{\omega_2}^T$$

and subject to these conditions, we now maximize  $E_{\theta}[\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}]$  for each  $\mathbf{t}$ , where  $\theta \neq \theta_0$ .

This conditional problem has already been solved in [Section 6.7.1](#), since for each  $\mathbf{t}$ ,  $\{f_{U|\mathbf{t}}(u; \theta) = C_{\mathbf{t}}(\theta) \exp[\theta u] h_{\mathbf{t}}(u), \theta \in I\}$  is an exponential family. The maximizing test is given as

$$\varphi_2(u, \mathbf{t}) = \begin{cases} 0 & \text{if } c_{21}(\mathbf{t}) < u < c_{22}(\mathbf{t}) \\ \gamma_{2i}(\mathbf{t}) & \text{if } u = c_{2i}(\mathbf{t}), i = 1, 2 \\ 1 & \text{if } u < c_{21}(\mathbf{t}) \text{ or } u > c_{22}(\mathbf{t}), \end{cases}$$

where  $c_{2i}(\mathbf{t})$  and  $0 \leq \gamma_{2i}(\mathbf{t}) \leq 1$ ,  $i = 1, 2$ , for each  $\mathbf{t}$  are determined by

$$E_{\theta_0}[\varphi_2(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha \text{ and } E_{\theta_0}[U\varphi_2(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha E_{\theta_0}[U | \mathbf{T} = \mathbf{t}].$$

To show that  $\varphi_2$  is a UMP level  $\alpha$  test for  $H_{20}$  vs  $H_{21}$ , we only need to verify that  $\varphi_2$  is a level  $\alpha$  test, as shown below:

$$\begin{aligned} E_{\theta_0}[\varphi_2(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] &= \alpha, \text{ a.e., } \mathcal{P}_{\omega_2}^T \\ \Rightarrow E_{\theta_0, \tau}[\varphi_2(U, \mathbf{T})] &= E_{\theta_0, \tau} E_{\theta_0}[\varphi_2(U, \mathbf{T}) | \mathbf{T}] = E_{\theta_0, \tau}[\alpha] = \alpha \quad \text{for all } \tau \in \omega_2. \end{aligned}$$

**Problem 3.** UMP Unbiased Level  $\alpha$  Test for  $H_{30}: \theta_1 \leq \theta \leq \theta_2$  vs  $H_{31}: \theta \notin [\theta_1, \theta_2]$ .

Here we maximize  $E_{\theta}[\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}]$ ,  $\theta \notin [\theta_1, \theta_2]$ , subject to  $E_{\theta_i}[\varphi(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha$ ,  $i = 1, 2$  for each  $\mathbf{t}$ . As seen in [Section 6.7.2](#), this results in the test:

$$\varphi_3(u, \mathbf{t}) = \begin{cases} 0 & \text{if } c_{31}(\mathbf{t}) < u < c_{32}(\mathbf{t}) \\ \gamma_{3i}(\mathbf{t}) & \text{if } u = c_{3i}(\mathbf{t}), i = 1, 2 \\ 1 & \text{if } u < c_{31}(\mathbf{t}) \text{ or } u > c_{32}(\mathbf{t}), \end{cases}$$

where  $c_{3i}(\mathbf{t})$  and  $0 \leq \gamma_{3i}(\mathbf{t}) \leq 1$ ,  $i = 1, 2$  for each  $\mathbf{t}$  is determined by

$$E_{\theta_i}[\varphi_3(U, \mathbf{t}) | \mathbf{T} = \mathbf{t}] = \alpha, \quad i = 1, 2.$$

## Examples of Conditional Tests

**Example 6.9.3.** Let  $\mathbf{X} = (X_1, \dots, X_m)$  and  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be independent random samples from  $Poisson(\lambda)$  and  $Poisson(\mu)$ , respectively. We want to test  $H_0: \mu = a\lambda$  vs  $H_1: \mu \neq a\lambda$  where  $a > 0$  is given. Here  $(X, Y) = (\sum_{i=1}^m X_i, \sum_{i=1}^n Y_i)$  is sufficient for  $(\lambda, \mu)$  in  $(\mathbf{X}, \mathbf{Y})$ . As in

**Example 6.9.2,** reparametrize and transform the data as:

$$\pi = n\mu/(m\lambda + n\mu), \xi = m\lambda + n\mu, \text{ and } U = Y, \quad T = X + Y,$$

to arrive at

$$f_{U,T}(u, t; \pi, \xi) = \binom{t}{u} \pi^u (1 - \pi)^{t-u} e^{-\xi} \xi^t / t!,$$

$t = 0, 1, \dots$  and  $u = 0, 1, \dots, t$ .

To put the problem in the framework of a two-parameter exponential family, we further reparametrize:

$$\theta = \log\left(\frac{\pi}{1 - \pi}\right) = \log\left(\frac{n\mu}{m\lambda}\right) \text{ and } \tau = \log(\xi(1 - \pi)) = \log(m\lambda)$$

to write the joint pmf of  $(U, T)$  as:

$$f_{U,T}(u, t; \theta, \tau) = \exp[-(e^\tau + e^{\theta+\tau})] e^{\theta u + \tau t} \frac{1}{u!(t-u)!}.$$

The hypotheses  $H_0$  and  $H_1$  now become  $H_0: \theta = \theta_0$  and  $H_1: \theta \neq \theta_0$ , where  $\theta_0 = \log(na/m)$ .

The UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  ( $\tau$  being a nuisance parameter) is

$$\varphi(u, t) = \begin{cases} 0 & \text{if } c_1(t) < u < c_2(t) \\ \gamma_i(t) & \text{if } u = c_i(t), \quad i = 1, 2 \\ 1 & \text{if } u < c_1(t) \text{ or } u > c_2(t), \end{cases}$$

where  $c_i(t)$  and  $0 \leq \gamma_i(t) \leq 1$ ,  $i = 1, 2$ , are determined by:

$$E_{\theta_0}[\varphi(U, t)|T = t] = \alpha \text{ and } E_{\theta_0}[U\varphi(U, t)|T = t] = \alpha E_{\theta_0}[U|T = t].$$

These equations can be written explicitly, using

$$f_{U|T}(u; \theta_0) = \binom{t}{u} \left(\frac{na}{m+na}\right)^u \left(\frac{m}{m+na}\right)^{t-u}, \quad u = 0, 1, \dots, t, \text{ and}$$

$$E_{\theta_0}[U|T = t] = t \left(\frac{na}{m+na}\right).$$

The problems of finding unbiased level  $\alpha$  tests for testing  $H_0: \mu \leq a\lambda$  vs  $H_1: \mu > a\lambda$  and for testing  $H_0: a_1\lambda \leq \mu \leq a_2\lambda$  vs  $H_1: \mu \notin [a_1\lambda, a_2\lambda]$  for given  $a$  or given  $a_1 < a_2$  are treated analogously using the methods developed for [Problems 1](#) and [3](#) of [Section 6.9](#).

**Example 6.9.4.** Let  $X \sim \text{Binomial}(m, \pi_1)$  and  $Y \sim \text{Binomial}(n, \pi_2)$  be independent. We want UMP unbiased level  $\alpha$  test for  $H_0: \pi_1 = \pi_2$  vs  $H_1: \pi_1 \neq \pi_2$ . Here,

$$f_{X,Y}(x,y;\pi_1,\pi_2) = \binom{m}{x} \binom{n}{y} (1-\pi_1)^m (1-\pi_2)^n \\ \times \exp \left[ \left( \log \frac{\pi_2(1-\pi_1)}{\pi_1(1-\pi_2)} \right) y + \left( \log \frac{\pi_1}{1-\pi_1} \right) (x+y) \right].$$

Reparametrizing

$$\theta = \log \frac{\pi_2(1-\pi_1)}{\pi_1(1-\pi_2)}, \quad \tau = \log \frac{\pi_1}{1-\pi_1}$$

and transforming the data to  $U = Y$ ,  $T = X + Y$ , the model is equivalently expressed as:

$$f_{U,T}(u,t;\theta,\tau) = C(\theta,\tau) \exp[\theta u + \tau t] \quad \text{and} \quad \pi_1 = \pi_2 \Leftrightarrow \theta = 0.$$

The UMP unbiased level  $\alpha$  test for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  is of the form as in [Example 6.9.3](#) with

$$f_{U|t}(u;\theta=0) = P_{\pi_1=\pi_2}[Y=u|X+Y=t] = \frac{\binom{m}{t-u}\binom{n}{u}}{\binom{m+n}{t}}$$

which is hypergeometric, and  $E_{\theta=0}[U|T=t] = \binom{m}{m+n}t$ .

**Example 6.9.5.**  $A$  and  $B$  are events in a probability space. The frequencies of  $AB$ ,  $A^cB$ ,  $AB^c$ , and  $A^cB^c$  in  $n$  trials are given in [Table 6.1](#), known as a  $2 \times 2$  contingency table:

**Table 6.1** Frequency Distribution in  $n$  Trials

Events	$A$	$A^c$	Subtotal
$B$	$X$	$Z$	$T$
$B^c$	$Y$	$W$	$T'$
Subtotal	$S$	$S'$	$n$

Based on this data, we want to test:  $H_0: A$  and  $B$  are *independent* or *negatively dependent* (ie,  $p_{AB} \leq p_A p_B$ ) vs  $H_1: A$  and  $B$  are *positively dependent* (ie,  $p_{AB} > p_A p_B$ ). Let  $\Delta = (p_{AB^c} p_{A^c B}) / (p_{AB} p_{A^c B^c})$ . Then the problem can be equivalently described as that of testing  $H_0: \Delta \geq 1$  vs  $H_1: \Delta < 1$ . Reparametrize by transforming

$$(p_{AB}, p_{AB^c}, p_{A^c B}) \rightarrow (\tau_0, \tau_1, \tau_2) \rightarrow (\theta, \tau_1, \tau_2),$$

where

$$\begin{aligned} \tau_0 &= \log(p_{AB}/p_{A^c B^c}), & \tau_1 &= \log(p_{AB^c}/p_{A^c B^c}), \\ \tau_2 &= \log(p_{A^c B}/p_{A^c B^c}), \text{ and } \theta = \tau_0 - \tau_1 - \tau_2 = -\log \Delta. \end{aligned}$$

Also transform the data from  $(X, Y, Z) \rightarrow (X, X+Y, X+Z) = (X, S, T)$ . In terms of  $(\theta, \tau_1, \tau_2)$ , we can write

$$\begin{aligned} p_{A^c B^c} &= (1 + e^{\theta+\tau_1+\tau_2} + e^{\tau_1} + e^{\tau_2})^{-1}, & p_{AB} &= e^{\theta+\tau_1+\tau_2} p_{A^c B^c}, \\ p_{AB^c} &= e^{\tau_1} p_{A^c B^c}, \quad \text{and} \quad p_{AB} &= e^{\tau_2} p_{A^c B^c}. \end{aligned}$$

The multinomial pmf of  $(X, Y, Z)$  can now be rewritten, using  $Y = S - X, Z = T - X$  and the above formulas for  $p_{AB}$ ,  $p_{AB^c}$ ,  $p_{A^cB}$ , and  $p_{A^cB^c}$  to express the pmf of  $(X, S, T)$  as:

$$f_{X,S,T}(x, s, t; \theta, \tau_1, \tau_2) = C(\theta, \tau_1, \tau_2) \exp[\theta x + \tau_1 s + \tau_2 t] h(x, s, t),$$

where  $C(\theta, \tau_1, \tau_2) = [1 + e^{\theta+\tau_1+\tau_2} + e^{\tau_1} + e^{\tau_2}]$  and  $h(x, s, t)$  is the multinomial coefficient

$$n!/\{x!(s-x)!(t-x)!(n-s-t+x)!\}.$$

Now the UMP unbiased level  $\alpha$  test for  $H_0: \Delta \geq 1$  vs  $H_1: \Delta < 1$  (ie, for  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$ ) is given by

$$\varphi(x, s, t) = \begin{cases} 0 & \text{if } x < c(s, t) \\ \gamma(s, t) & \text{if } x = c(s, t) \\ 1 & \text{if } x > c(s, t), \end{cases}$$

where  $c(s, t)$  and  $0 \leq \gamma(s, t) \leq 1$  are determined by

$$\gamma(s, t)P_{\theta=0}[X = c(s, t)|S = s, T = t] + P_{\theta=0}[X > c(s, t)|S = s, T = t] = \alpha.$$

Note that  $P_{\theta=0}[X = x|S = s, T = t]$

$$= \begin{cases} \binom{t}{x} \binom{n-t}{s-x} / \binom{n}{s}, & \max(0, s+t-n) \leq x \leq \min(n, s) \\ 0 & \text{otherwise} \end{cases}$$

is the hypergeometric distribution.

The UMP unbiased level  $\alpha$  test for  $H_0: p_{AB} = p_A p_B$  (independence) vs  $H_1: p_{AB} \neq p_A p_B$  (ie,  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$ ) is obtained by the same approach. This is known as the Fisher-Irwin test (also called the “Fisher exact test”), which is formally the same as the test obtained in [Example 6.9.2](#).

## Simplified Versions of Conditional Tests

The conditional tests  $\varphi_1(u, t)$  for  $H_{10}: \theta \leq \theta_0$  vs  $H_{11}: \theta > \theta_0$ ,  $\varphi_2(u, t)$  for  $H_{20}: \theta = \theta_0$  vs  $H_{21}: \theta \neq \theta_0$  and  $\varphi_3(u, t)$  for  $H_{30}: \theta_1 \leq \theta \leq \theta_2$  vs  $H_{31}: \theta \notin [\theta_1, \theta_2]$  are inconvenient in many applications. However, in some situations which include testing problems in the family of normal distributions, these tests can be equivalently expressed in terms of a single statistic  $V = g(U, T)$ , as shown in the following theorem.

**Theorem 6.9.2.**

**A.** Suppose  $V = g(U, T)$  is such that

- (a)  $V$  is independent of  $T$  when  $\theta = \theta_0$ , and
- (b)  $g(u, t)$  is increasing in  $u$  for each  $t$ .

Then the UMP unbiased level  $\alpha$  test  $\varphi_1$  for [Problem 1](#) becomes

$$\varphi_1^*(v) = \begin{cases} 0 & \text{if } v < c \\ \gamma & \text{if } v = c \\ 1 & \text{if } v > c, \end{cases}$$

where  $c$  and  $0 \leq \gamma \leq 1$  are determined by  $E_{\theta_0}[\varphi_1^*(V)] = \alpha$ .

**B.** Suppose  $V = g(U, T)$  is such that

- (a)  $V$  is independent of  $T$  when  $\theta = \theta_0$ , and
- (b)  $g(u, t) = a(t)u + b(t)$  with  $a(t) > 0$ .

Then the UMP unbiased level  $\alpha$  test  $\varphi_2$  for [Problem 2](#) becomes

$$\varphi_2^*(v) = \begin{cases} 0 & \text{if } c_1 < v < c_2 \\ \gamma_i & \text{if } v = c_i, \quad i = 1, 2 \\ 1 & \text{if } v < c_1 \text{ or } v > c_2, \end{cases}$$

where  $c_i$  and  $0 \leq \gamma_i \leq 1$ ,  $i = 1, 2$  are determined by  $E_{\theta_0}[\varphi_2^*(V)] = \alpha$  and  $E_{\theta_0}[V\varphi_2^*(V)] = \alpha E_{\theta_0}[V]$ .

**C.** Suppose  $V = g(U, T)$  is such that

- (a)  $V$  is independent of  $T$  when  $\theta = \theta_1$  and when  $\theta = \theta_2$ , and
- (b)  $g(u, t)$  is increasing in  $u$  for each  $t$ .

Then the UMP unbiased level  $\alpha$  test  $\varphi_3$  in [Problem 3](#) becomes  $\varphi_3^*$  having the same form as  $\varphi_2^*$  (in part B of the theorem), but  $c_i$  and  $0 \leq \gamma_i \leq 1$ ,  $i = 1, 2$  are determined by  $E_{\theta_i}[\varphi_3^*(V)] = \alpha$ ,  $i = 1, 2$ .

*Proof of Part A.* The UMP unbiased level  $\alpha$  test  $\varphi_1$  in [Problem 1](#) can be equivalently expressed in the form of  $\varphi_1^*$  by condition (b), but the quantities  $c$  and  $0 \leq \gamma \leq 1$  may depend on  $t$  and are determined by

$$\gamma(t)P_{\theta_0}[V = c(t)|T = t] + P_{\theta_0}[V > c(t)|T = t] = \alpha \quad \text{for all } t.$$

By condition (a),  $V$  is independent of  $T$  when  $\theta = \theta_0$ , so  $c(t)$  and  $0 \leq \gamma(t) \leq 1$  do not depend on  $t$  and are determined by  $E_{\theta_0}[\varphi_1^*(V)] = \alpha$ .  $\square$

The proof of part C is analogous.

*Proof of Part B.* Here also, by condition (b), the UMP unbiased level  $\alpha$  test  $\varphi_2$  in [Problem 2](#) can be equivalently expressed in the form  $\varphi_2^*$ , but  $c_i(t)$  and  $0 \leq \gamma_i \leq 1$ ,  $i = 1, 2$  may depend on  $t$  and are determined by

- (i)  $E_{\theta_0}[\varphi_2^*(V, t)|T = t] = \alpha$  and
- (ii)  $E_{\theta_0}[U\varphi_2^*(V, t)|T = t] = \alpha E_{\theta_0}[U|T = t]$  for all  $t$ .

Substituting  $V = a(t)U + b(t)$  (ie,  $U = (V - b(t))/a(t)$  in (ii)), we have

$$\begin{aligned} a(t)^{-1}E_{\theta_0}[\{V - b(t)\}\varphi_2^*(V, t)|T = t] &= \alpha a(t)^{-1}E_{\theta_0}[V - b(t)|T = t] \quad \text{for all } t \\ \Leftrightarrow E_{\theta_0}[V\varphi_2^*(V, t)|T = t] - b(t)E_{\theta_0}[\varphi_2^*(V, t)|T = t] &= \alpha E_{\theta_0}[V|T = t] - \alpha b(t) \quad \text{for all } t \\ \Leftrightarrow E_{\theta_0}[V\varphi_2^*(V, t)|T = t] &= \alpha E_{\theta_0}[V|T = t] \quad \text{for all } t, \end{aligned}$$

since  $E_{\theta_0}[\varphi_2^*(V, t)|T = t] = \alpha$  for all  $t$  by (i). Now using condition (a), we see that  $c_i(t)$  and  $0 \leq \gamma_i(t) \leq 1$ ,  $i = 1, 2$  do not depend on  $t$  and are determined by:  $E_{\theta_0}[\varphi_2^*(V)] = \alpha$  and  $E_{\theta_0}[V\varphi_2^*(V)] = \alpha E_{\theta_0}[V]$ .  $\square$

In actual applications, condition (b) of this theorem is verified in an obvious manner. The verification of condition (a) often follows from  $V$  being an ancillary statistic (ie, one whose distribution does not depend on the nuisance parameter  $\tau$ ) and  $T$  being a complete sufficient statistic for  $\tau$ , using Basu's [Theorem 5.1.6](#).

## Examples: Hypothesis Testing With Nuisance Parameters in the Context of Normal Distribution

In each of the following examples, we go through the following steps:

- I. The data which is a random sample from the population(s) under consideration are:  $(X_1, \dots, X_n)$ , or  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  or  $((X_1, Y_1), \dots, (X_n, Y_n))$ . Also state  $H_0$  and  $H_1$  in terms of the population parameters.
- II. The sufficient statistics are  $(U, T)$ . The conditional distribution of  $U$  given  $T = t$  is complicated, so the conditional test described earlier in [Section 6.9](#) are impractical.
- III. Reparametrize from  $(\mu, \sigma^2)$  or  $(\mu_1, \mu_2, \sigma^2)$ , etc., to  $(\theta, \tau)$  and restate  $H_0$  and  $H_1$  in terms of  $\theta$ , treating  $\tau$  as the nuisance parameter. In each example, the pdf/pmf of  $(U, T)$  belongs to the exponential family with natural parameters  $(\theta, \tau)$ .
- IV. Transform  $(U, T)$  to  $V = g(U, T)$  or  $(V = g_1(U, T), W = g_2(U, T))$  as needed to use [Theorem 6.9.2](#), verifying the conditions.
- V. Express the UMP unbiased level  $\alpha$  test in terms of  $V$  or  $(V, W)$  using the distributions of  $V$  and  $W$  under  $H_0$ .

### Example 6.9.6.

Testing for  $\mu$  of  $N(\mu, \sigma^2)$  with  $\sigma^2$  unknown.

- I. On the basis of the data  $(X_1, \dots, X_n)$  test (i)  $H_0: \mu \leq \mu_0$  vs  $H_1: \mu > \mu_0$  and (ii)  $H_0: \mu = \mu_0$  vs  $H_1: \mu \neq \mu_0$  where  $\mu_0 = 0$  (otherwise, replace  $\mu$  by  $\mu - \mu_0$ ).
- II. The sufficient statistics are  $(U, T) = (\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ .
- III. Reparametrize  $\theta = \mu/\sigma^2$ ,  $\tau = -1/(2\sigma^2)$ . In (i), test  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  where  $\theta_0 = \mu_0/\sigma^2$  and in (ii), test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ ,  $\tau$  being the nuisance parameter.
- IV. Express the UMP unbiased level  $\alpha$  tests  $\varphi_1(U, T)$  for (i) as  $\varphi_1^*(V)$  and  $\varphi_2(U, T)$  for (ii) as  $\varphi_2^*(W)$ , where

$$V = g_1(U, T) = \sqrt{\frac{n-1}{n}} \frac{U}{\sqrt{T - U^2/n}} = \frac{\sqrt{n}\bar{X}}{s} \text{ and } W = g_2(U, T) = U/\sqrt{T}.$$

Verify conditions (a) and (b) of [Theorem 6.9.2A](#) for  $V$  and  $W$ .

- V. The test for (i) is  $\varphi_1^*(v) = I_{[c, \infty)}(v)$ , where  $c$  satisfies  $P_{\mu=0}[\sqrt{n}\bar{X}/s \geq c] = \alpha$ . Under  $\mu = 0$ ,  $\sqrt{n}\bar{X}/s \sim t_{n-1}$ , so  $c = t_{n-1}(\alpha)$  where  $P[t_{n-1} \geq t_{n-1}(\alpha)] = \alpha$ . The test for (ii) is  $\varphi_2^*(w) = I_{[c, \infty)}(|w|)$ , using the fact that under  $\mu = 0$ ,  $W$  is symmetrically distributed about 0, and  $c$  satisfies  $P_{\mu=0}[|W| \geq c] = \alpha$ . Finally, note that

$$V = \sqrt{\frac{n-1}{n}} W / \sqrt{1 - W^2/n} = \sqrt{n}\bar{X}/s \sim t_{n-1} \text{ when } \mu = 0$$

and  $|V|$  is an increasing function of  $|W|$ . Thus  $\varphi_2^*$  becomes  $\varphi_2^*(v) = I_{[t_{n-1}(\alpha/2), \infty)}(|v|)$ .

### Example 6.9.7. Testing for $\sigma^2$ in $N(\mu, \sigma^2)$ with $\mu$ unknown.

- I. On the basis of the data  $(X_1, \dots, X_n)$  test (i)  $H_0: \sigma^2 \leq \sigma_0^2$  vs  $H_1: \sigma^2 > \sigma_0^2$  and (ii)  $H_0: \sigma^2 = \sigma_0^2$  vs  $H_1: \sigma^2 \neq \sigma_0^2$ .
- II. Here  $(U, T) = (\sum_1^n X_i^2, \sum_1^n X_i)$  are the sufficient statistics.

- III. Reparametrize  $\theta = -1/(2\sigma^2)$  and  $\tau = \mu/\sigma^2$ . In (i), test  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  where  $\theta_0 = -1/(2\sigma_0^2)$  and in (ii), test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .
- IV. Express the UMP unbiased level  $\alpha$  tests  $\varphi_1(U, T)$  for (i) as  $\varphi_1^*(V)$  and  $\varphi_2(U, T)$  for (ii) as  $\varphi_2^*(V)$ , where

$$V = g(U, T) = (1/\sigma_0^2)(U - T^2/n) = \sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma_0} \right)^2.$$

Verify conditions (a) and (b) of [Theorem 6.9.2A](#) for  $V$ .

- V. Since  $V \sim \chi_{n-1}^2$  under  $\sigma^2 = \sigma_0^2$ , the UMP level  $\alpha$  test for (i) is  $\varphi_1^*(v) = I_{[c, \infty)}(v)$  where  $c = \chi_{n-1}^2(\alpha)$ ,  $P[\chi_{n-1}^2 \geq \chi_{n-1}^2(\alpha)] = \alpha$ , and for (ii),  $\varphi_2^*(v) = I_{(0, c_1]}(v) + I_{[c_2, \infty)}(v)$  where  $c_1 < c_2$  satisfy

$$\int_{c_1}^{c_2} f_{n-1}(y) dy = 1 - \alpha \text{ and } \int_{c_1}^{c_2} y f_{n-1}(y) dy = (n-1)(1 - \alpha),$$

$f_{n-1}$  being the pdf of  $\chi_{n-1}^2$ . The second condition in the display above can be expressed in any of the following two ways (as in [Example 6.7.2](#)):

$$\begin{aligned} \int_{c_1}^{c_2} f_{n+1}(y) dy &= 1 - \alpha \text{ since } y f_{n-1}(y) = (n-1)f_{n+1}(y), \text{ or} \\ e^{-c_1/2} c_1^{(n-1)/2} &= e^{-c_2/2} c_2^{(n-1)/2} \text{ since} \\ \int_{c_1}^{c_2} y f_{n-1}(y) dy &= \frac{1}{2^{(n-3)/2} \Gamma((n-1)/2)} \left( e^{-c_1/2} c_1^{(n-1)/2} - e^{-c_2/2} c_2^{(n-1)/2} \right) \\ &\quad + n(n-1)(1 - \alpha), \end{aligned}$$

using integration by parts and condition (i). For moderately large  $n$  and  $\theta_0$  neither too small nor too large, CLT provides a reasonably good approximation for  $c_1$  and  $c_2$  determined by

$$\int_0^{c_1} f_{n-1}(y) dy = \int_{c_2}^{\infty} f_{n-1}(y) dy = \alpha/2.$$

**Example 6.9.8.** Testing for  $\mu_1 = \mu_2$  of  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$  with  $\mu_1, \mu_2, \sigma^2$  unknown.

- I. On the basis of the data  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  from  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$ , respectively, test
  - (i)  $H_0: \mu_2 - \mu_1 \leq 0$  vs  $H_1: \mu_2 - \mu_1 > 0$  and
  - (ii)  $H_0: \mu_2 - \mu_1 = 0$  vs  $H_1: \mu_2 - \mu_1 \neq 0$ .
- II. The sufficient statistics and their one-to-one transforms are

$$\begin{aligned} &\left( \sum_1^n Y_i, \sum_1^m X_i, \sum_1^m X_i^2 + \sum_1^n Y_i^2 \right) \text{ and} \\ &(U, T_1, T_2) = \left( \bar{Y} - \bar{X}, m\bar{X} + n\bar{Y}, \sum_1^m X_i^2 + \sum_1^n Y_i^2 \right) \end{aligned}$$

### III. Reparametrize

$$(\theta, \tau_1, \tau_2) = \left( \frac{\mu_2 - \mu_1}{(m^{-1} + n^{-1})\sigma^2}, \frac{m\mu_1 + n\mu_2}{(m+n)\sigma^2}, -\frac{1}{2\sigma^2} \right).$$

In (i), test  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$  where  $\theta_0 = 0$  and in (ii), test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ ,  $\tau_1, \tau_2$  being the nuisance parameters.

**IV, V.** For testing  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$ , that is,  $H_0: \mu_2 \leq \mu_1$  vs  $H_1: \mu_2 > \mu_1$ , take

$$\begin{aligned} V &= \sqrt{\frac{mn(m+n-2)}{m+n}} \frac{U}{\sqrt{T_2 - (T_1^2 + mnU^2)/(m+n)}} \\ &= \frac{\bar{Y} - \bar{X}}{\sqrt{\frac{1}{m} + \frac{1}{n}} \sqrt{\frac{\sum_{i=1}^m (X_i - \bar{X})^2 + \sum_{i=1}^n (Y_i - \bar{Y})^2}{m+n-2}}} := \frac{\bar{Y} - \bar{X}}{\sqrt{\frac{1}{m} + \frac{1}{n}}} \dots \end{aligned}$$

Then  $V$  satisfies both conditions of [Theorem 6.9.2A](#), and the UMP unbiased level  $\alpha$  test for  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  is  $\varphi_1^*(v) = I_{[c, \infty)}(v)$ , where  $c = t_{m+n-2}(\alpha)$ .  
For testing  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$ , take

$$\begin{aligned} W &= \frac{U}{\sqrt{T_2 - T_1^2/(m+n)}} \\ &= \frac{\bar{Y} - \bar{X}}{\sqrt{\sum_1^m X_i^2 + \sum_1^n Y_i^2 - \{\sum_1^m X_i + \sum_1^n Y_i\}^2/(m+n)}}. \end{aligned}$$

Then  $W$  satisfies both conditions of [Theorem 6.9.2B](#), and the UMP unbiased level  $\alpha$  test for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  is  $\varphi_2^*(w) = I_{[c, \infty)}(|w|)$ , where  $c$  satisfies  $P_{\mu_1=\mu_2}[|W| \geq c] = \alpha$ . Now let  $V = W/\sqrt{1 - mnW^2/(m+n)}$ , which is the same as  $V$  defined for problem (i) and  $|V|$  is an increasing function of  $|W|$ . Thus the UMP unbiased level  $\alpha$  test for  $H_0: \mu_1 = \mu_2$  vs  $H_1: \mu_1 \neq \mu_2$  becomes  $\varphi_2^*(v) = I_{[c, \infty)}(|v|)$ , where  $c = t_{m+n-2}(\alpha/2)$ .

**Example 6.9.9.** Testing for  $\sigma_1^2 = \sigma_2^2$  of  $N(\mu_1, \sigma_1^2)$  and  $N(\mu_2, \sigma_2^2)$ , all parameters are unknown.

- I. On the basis of the data  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  from  $N(\mu_1, \sigma_1^2)$  and  $N(\mu_2, \sigma_2^2)$ , respectively, test (for a given  $k > 0$ ):

- (i)  $H_0: \sigma_2^2 \leq k\sigma_1^2$  vs  $H_1: \sigma_2^2 > k\sigma_1^2$  and  
(ii)  $H_0: \sigma_2^2 = k\sigma_1^2$  vs  $H_1: \sigma_2^2 \neq k\sigma_1^2$ .

- II. The sufficient statistics and their one-to-one transforms are

$$\begin{aligned} &\left( \sum_1^n Y_i^2, \sum_1^m X_i^2, \sum_1^n Y_i, \sum_1^m X_i \right) \text{ and} \\ &(U, T_1, T_2, T_3) = \left( \sum_1^n Y_i^2, \sum_1^m X_i^2 + k^{-1} \sum_1^n Y_i^2, \sum_1^n Y_i, \sum_1^m X_i \right). \end{aligned}$$

### III. Reparametrize

$$(\theta, \tau_1, \tau_2, \tau_3) = \left( -\frac{1}{2\sigma_2^2} + \frac{1}{2k\sigma_1^2}, -\frac{1}{2\sigma_1^2}, \frac{\mu_2}{\sigma_2^2}, \frac{\mu_1}{\sigma_1^2} \right),$$

$\tau_1, \tau_2, \tau_3$  being the nuisance parameters. In (i), test  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  and in (ii), test  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$ .

### IV, V. Take

$$V = \frac{\{(n-1)k\}^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2}{(m-1)^{-1} \sum_{i=1}^m (X_i - \bar{X})^2}$$

in problem (i). Then both conditions of [Theorem 6.9.2A](#) are satisfied and the UMP unbiased level  $\alpha$  test in problem (i) is  $\varphi_1^*(v) = I_{[c, \infty)}(v)$ , where  $P[F_{n-1, m-1} \geq c] = \alpha$ . In problem (ii), for testing  $H_0: \sigma_2^2 = k\sigma_1^2$  vs  $H_1: \sigma_2^2 \neq k\sigma_1^2$ , proceed with

$$W = k^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 / \left[ \sum_{i=1}^m (X_i - \bar{X})^2 + k^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \right]$$

to set up the test and then transform to  $V$  which increases with  $W$ .

**Example 6.9.10.** Testing for the slope parameter in Simple Linear Regression Model.

- I. On the basis of the data  $((x_1, Y_1), \dots, (x_n, Y_n))$  where  $x_1, \dots, x_n$  are given constants and  $Y_i = \alpha + \beta x_i + \epsilon_i$ ,  $\epsilon_1, \dots, \epsilon_n$  being iid  $N(0, \sigma^2)$ , with  $\alpha, \beta, \sigma^2$  all unknown, we want to test

- (i)  $H_0: \beta \leq \beta_0$  vs  $H_1: \beta > \beta_0$  and
- (ii)  $H_0: \beta = \beta_0$  vs  $H_1: \beta \neq \beta_0$ .

Without loss of generality, let  $\beta_0 = 0$  (otherwise, replace  $Y_i$  by  $Y_i - \beta_0 x_i$ ).

1. Transform the data to reduce the model to its canonical form: Let  $Y_i^* = Y_i - \alpha - \beta x_i$ ,  $i = 1, \dots, n$  and let  $\mathbf{a}_i^\top = (a_{i1}, \dots, a_{in})$ ,  $i = 1, \dots, n$  be  $n$ -dim vectors which form an orthonormal basis for  $\mathbb{R}^n$  (ie,  $\mathbf{a}_i^\top \mathbf{a}_i = 1$  for all  $i$  and  $\mathbf{a}_i^\top \mathbf{a}_j = 0$  for all  $i \neq j$ ). In particular, choose the first two vectors as

$$\mathbf{a}_1^\top = n^{-1/2}(1, \dots, 1) \text{ and } \mathbf{a}_2^\top = S_{xx}^{-1/2}(x_1 - \bar{x}, \dots, x_n - \bar{x}),$$

where  $\bar{x} = n^{-1} \sum_1^n x_i$  and  $S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$ . Let  $\mathbf{Y}^* = (Y_1^*, \dots, Y_n^*)$ , and define

$$W_1 = \mathbf{a}_1^\top \mathbf{Y}^*, W_2 = \mathbf{a}_2^\top \mathbf{Y}^* \text{ and } W_3 = \sum_{i=1}^n a_{3i} Y_i^*, \dots, W_n = \sum_{i=1}^n a_{ni} Y_i^*.$$

Then  $W_1, \dots, W_n$  are iid  $N(0, \sigma^2)$  and  $\sum_1^n W_i^2 = \sum_1^n Y_i^{*2}$ . Define  $S_{xY} = \sum_{i=1}^n (x_i - \bar{x}) Y_i$  and  $S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2$ . It follows that

- (a)  $\bar{Y} = n^{-1} \sum_1^n Y_i \sim N(\alpha + \beta \bar{x}, \sigma^2/n)$  because  $W_1 \sim N(0, \sigma^2)$ ,
- (b)  $W_2 \sim N(0, \sigma^2)$  (ie,  $B = S_{xY}/S_{xx} \sim N(\beta, \sigma^2/S_{xx})$ ),

- (c)  $R = \sum_{i=3}^n W_i^2 = \sum_{i=1}^n Y_i^{*2} - W_1^2 - W_2^2 = S_{YY} - S_{xY}^2/S_{xx} \sim \sigma^2 \chi_{n-2}^2$ ,  
(d)  $\bar{Y}$ ,  $B$ , and  $R$  are mutually independent.

**II.** The sufficient statistics and their one-to-one transforms are

$$(B, \bar{Y}, R) \text{ and } (U, T_1, T_2) = (B, \bar{Y}, n\bar{Y}^2 + S_{xx}B^2 + R).$$

**III.** Reparametrize

$$\theta = S_{xx}\beta/\sigma^2, \tau_1 = n(\alpha + \beta\bar{x})/\sigma^2, \tau_2 = -1/(2\sigma^2).$$

Starting from the joint distribution of  $(\bar{Y}, B, R)$ , using (a)-(d) in Step I, the joint pdf of  $(U, T_1, T_2)$  can be seen to be of the form

$$f_{U, T_1, T_2}(u, t_1, t_2; \theta, \tau_1, \tau_2) = C(\theta, \tau_1, \tau_2) \exp[\theta u + \tau_1 t_1 + \tau_2 t_2] h(u, t_1, t_2),$$

and the null and the alternative hypotheses can be restated as  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  in problem (i) and  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  in problem (ii).

**IV, V.** In problem (i), take

$$V = g(U, T_1, T_2) = \frac{\sqrt{n-2}U}{\sqrt{(T_2 - nT_1^2)/S_{xx} - U^2}} = \frac{\sqrt{n-2}B}{\sqrt{R/S_{xx}}}.$$

Then  $V$  satisfies the conditions of [Theorem 6.9.2A](#) and it is distributed as  $t_{n-2}$  under  $H_0$ . Hence the UMP unbiased level  $\alpha$  test for  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  is

$$\varphi_1^*(v) = I_{[t_{n-2}(\alpha), \infty)}(v).$$

In problem (ii), first take

$$V_1 = g(U, T_1, T_2) = U \sqrt{\frac{(T_2 - nT_1^2)/S_{xx}}{B^2 + R/S_{xx}}} = B \sqrt{\frac{B^2 + R/S_{xx}}{1 + B^2/S_{xx}}}.$$

Then  $V_1$  satisfies the conditions of [Theorem 6.9.2B](#). Hence the UMP unbiased level  $\alpha$  test for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  is  $\varphi(v_1) = I_{[c, \infty)}(|v_1|)$ , where  $c$  satisfies  $P_{\theta=0}[|V_1| \geq c] = \alpha$ . Finally, let

$$V = \frac{\sqrt{n-2}V_1}{\sqrt{1 - V_1^2}} = \frac{\sqrt{n-2}B}{\sqrt{R/S_{xx}}}$$

as in problem (i). Then  $V \sim t_{n-2}$  under  $\theta = 0$  and  $|V|$  is an increasing function of  $|V_1|$ , so  $\varphi(v_1)$  is the same as  $\varphi^*(v) = I_{[t_{n-2}(\alpha/2), \infty)}(|v|)$ .

*Remark 6.9.1.* If  $\beta_0 \neq 0$ , then  $S_{xY}$  and  $S_{YY}$  are calculated from  $(x_i, Y_i - \beta_0 x_i)$ ,  $i = 1, \dots, n$  and become

$$S_{xY}^* = S_{xY} - \beta_0 S_{xx} \text{ and } S_{YY}^* = S_{YY} + \beta_0^2 S_{xx} - 2\beta_0 S_{xY},$$

respectively, and  $B$  and  $R$  are replaced by

$$B^* = B - \beta_0 \text{ and } R^* = S_{YY}^* - (S_{xY}^*)^2/S_{xx} = S_{YY} - S_{xY}^*/S_{xx} = R.$$

As a result, the statistic  $V$  is replaced by  $V^* = \sqrt{n-2}(B - \beta_0)/\sqrt{R/S_{xx}}$ .

## 6.10 The $p$ -Value: Another Way to Report the Result of a Test

In Example 6.7.1(a) of Section 6.7.2, let  $H_0: \mu = \mu_0 = 10$  vs  $H_1: \mu \neq 10$  with known  $\sigma^2 = 4$ , and suppose that in a random sample of size  $n = 25$ , the observed sample mean is  $\bar{x} = 10.8$ . Then the test statistic  $T = \sqrt{n}(\bar{X} - \mu_0)$  has observed value

$$T(x) = \sqrt{n}(\bar{x} - \mu_0)/\sigma = 5(10.8 - 10)/2 = 2.0.$$

Whether to accept or reject  $H_0$ , now depends on the choice of  $\alpha$ . For  $\alpha = 0.05$ ,  $\Phi(1 - \alpha/2) = 1.96$  and for  $\alpha = 0.01$ ,  $\Phi(1 - \alpha/2) = 2.575$ ; so for the same data we should reject  $H_0$  at level  $\alpha = 0.05$  but accept  $H_0$  at level  $\alpha = 0.01$ . This raises two concerns:

- (i) The choice of  $\alpha$  is subjective, contrary to the aim at objectivity in the theory of hypothesis testing.
- (ii) Even for a choice of chosen  $\alpha$ , the test procedure only tells us whether to “accept  $H_0$ ” or “reject  $H_0$ ” without any indication of how strongly the data favors the decision.

To address these concerns, it would be desirable to report not only whether  $H_0$  is rejected or accepted at a preassigned level  $\alpha$  by the observed value of the test statistic  $T(x)$  but also the smallest level of significance at which  $T(x)$  would reject  $H_0$ , which is called the  $p$ -value of  $T(x)$ .

**Definition 6.10.1.** For a problem of testing  $H_0$  vs  $H_1$ , if  $T(x)$  is the appropriate test statistic based on the observed data  $x$ , then  $p$ -value of  $T$  = minimum  $\alpha$  for which  $H_0$  would be rejected by  $T(x)$ .

For a one-sided level  $\alpha$  test  $\phi(x) = I_{[c_\alpha, \infty)}(T(x))$ ,

$$p\text{-value of } T = \min_{\alpha} \{T(x) \geq c_\alpha\} = P_{H_0}[T \geq T(x)],$$

and for a symmetric two-sided test level  $\alpha$  test  $\phi(x) = 1 - I_{(-c_\alpha, c_\alpha)}(T(x))$ ,

$$p\text{-value of } T = \min_{\alpha} \{|T(x)| \geq c_\alpha\} = P_{H_0}[T \geq |T(x)|] + P_{H_0}[T \leq -|T(x)|].$$

In Example 6.7.1(a) discussed above,

$$T = \sqrt{n}(\bar{X} - \mu_0)/\sigma \quad \text{and} \quad T(x) = 2.0.$$

Since  $T \sim N(0, 1)$  under  $H_0$ ,

$$p\text{-value of } T = P_{H_0}[|T| \geq 2.0] = 2\{1 - \Phi(2)\} = (2)(0.0228) = 0.0456.$$

This statement is more informative than the one saying that  $H_0$  is accepted at level  $\alpha = 0.01$  but rejected at  $\alpha = 0.05$ .

### 6.10.1 Pearson's $P_\lambda$ Statistic

The  $p$ -value of a statistic  $T$  has so far been discussed in the context of  $T(x)$  for an observed sample  $x$ . We now look at the  $p$ -value as a random variable:

$$p\text{-value of } T = P_{H_0}[T \geq T(X)]$$

considering a one-sided test for convenience of discussion.

Let  $F_T$  denote the cdf of  $T$  under  $H_0$ . Then

$$P_{H_0}[T \geq T(X)] = 1 - F_T(T(X)) \sim 1 - U \stackrel{D}{=} U,$$

where  $U \sim \text{Uniform}(0, 1)$  as seen in [Section 1.11](#).

In problem 16(a) of [Chapter 2](#), we have seen that

$$W = -2 \log U \sim \text{Exp}(1/2) \stackrel{D}{=} \chi_2^2.$$

Thus the  $p$ -value of  $T$  for a random sample  $X$  being  $P = P_{H_0}[T \geq T(X)]$ ,  $W = -2 \log P$  is distributed as  $\chi_2^2$ .

Now suppose that a certain study, such as a trial on the effectiveness of a drug, is carried out independently by several investigators. These studies, being designed differently and based on different sample sizes, may have been analyzed by different test statistics aimed at testing the same null hypothesis  $H_0$ , the rejection of which would indicate effectiveness of the drug. Suppose that the results of these studies are summarized by test statistics  $T_1, \dots, T_k$  with observed values  $T_1(x_1), \dots, T_k(x_k)$ , and let the observed  $p$ -values of these tests (using right-tail tests for simplicity) be denoted by

$$P(x_j) = P_{H_0}[T_j \geq T_j(x_j)], \quad j = 1, \dots, k.$$

Then these  $p$ -values are sample realizations of

$$P_j = P_{H_0}[T_j \geq T_j(X_j)]$$

and therefore,  $w_j = -2 \log P(x_j)$  are sample realizations of  $W_j = -2 \log P_j$ , which are iid as  $\chi_2^2$  under  $H_0$ .

Now define  $P_\lambda = \sum_{j=1}^k [-2 \log P_j]$  as a test statistic combining the results of all the  $k$  tests, which is distributed as  $\chi_{2k}^2$  under  $H_0$ , and can be used as a test statistic based on the combined evidence provided by the results of all the  $k$  investigations.

## 6.11 Sequential Probability Ratio Test

Traditional statistical inference deals with analysis of a set of data  $(X_1, \dots, X_n)$  to draw some conclusion about a parameter  $\theta$  involved in the joint pdf (or pmf)  $f_n(x_1, \dots, x_n; \theta)$ . The observations  $X_1, \dots, X_n$  are often iid with individual pdf  $f(x; \theta)$  (although they need not be so), in which case,  $f_n(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$ . But the main thing is that  $n$  is fixed and the functional form  $f_n$  is fixed (ie, the sample size and the sampling design are predetermined). However, some samples provide more conclusive evidence than others, so it makes sense to take samples gradually (ie, sequentially), and examine the evidence gathered at each stage of sampling to determine whether any more samples are needed, and if not, to draw a final conclusion. This is known as *Sequential Analysis*.

Here we shall only discuss the problem of testing a simple null hypothesis  $H_0: \theta = \theta_0$  against a simple alternative  $H_1: \theta = \theta_1$  based on iid samples  $X_1, X_2, \dots$  observed sequentially. For this, we need a procedure of the following form:

**Step 1.** Start with observing  $X_1$ . Let  $S_{10}, S_{11}$ , and  $C_1$  be a partition of the sample space  $\mathcal{X}_1 = \mathcal{X}$  of  $X_1$  such that  $S_{10} \cup S_{11} \cup C_1 = \mathcal{X}_1$ . If the observed value  $x_1 \in S_{10}$ , stop sampling and accept  $H_0$ ; if  $x_1 \in S_{11}$ , stop sampling and accept  $H_1$ ; if  $x_1 \in C_1$  (continuation region), take another observation  $X_2$ .

**Step 2.** Now the observed data are  $(x_1, x_2) \in C_1 \times \mathcal{X} = \mathcal{X}_2$ . Let  $S_{20}, S_{21}$ , and  $C_2$  be a partition of  $\mathcal{X}_2$ . If  $(x_1, x_2) \in S_{20}$ , stop sampling and accept  $H_0$ ; if  $(x_1, x_2) \in S_{21}$ , stop sampling and accept  $H_1$ ; if  $(x_1, x_2) \in C_2$  (continuation region), take another observation  $X_3$ .

**Step  $n$ .** At the  $n$ th stage of sampling, the observed data are  $(x_1, \dots, x_{n-1}, x_n) \in C_{n-1} \times \mathcal{X} = \mathcal{X}_n$ . Let  $S_{n0}, S_{n1}$ , and  $C_n$  be a partition of  $\mathcal{X}_n$ . If  $(x_1, \dots, x_n) \in S_{n0}$ , stop sampling and accept  $H_0$ ; if  $(x_1, \dots, x_n) \in S_{n1}$ , stop sampling and accept  $H_1$ ; if  $(x_1, \dots, x_n) \in C_n$ , observe  $X_{n+1}$ .

Sequential Probability Ratio Tests, which we shall refer to as SPRT, are due to Wald [24]. The SPRT is based on the following idea. For a fixed sample size  $n$ , the MP test for  $H_0$  vs  $H_1$  at a given level  $\alpha$ , accepts  $H_0$  or  $H_1$  according as

$$\lambda_n = \frac{f_n(x_1, \dots, x_n; \theta_1)}{f_n(x_1, \dots, x_n; \theta_0)} \text{ is } < k \text{ or } > k,$$

where  $k = k_n(\alpha)$  depends on the sample size and the level of significance. [For given  $\alpha, \beta$ , we can also determine the smallest sample size  $N(\alpha, \beta)$  such that for  $n \geq N(\alpha, \beta)$ , the MP level  $\alpha$  test based on a sample size  $n$  will have the probability of Type II error  $\leq \beta$ ]. In the sequential setting, at the  $n$ th stage of sampling, we modify the MP test described above by accepting  $H_0$  if the likelihood ratio  $\lambda_n \leq B$ , accepting  $H_1$  if  $\lambda_n \geq A$ , and observing  $X_{n+1}$  if  $B < \lambda_n < A$ .

### Definition of $SPRT(A, B)$ , $B < 1 < A$ , Based on iid Observations

Let  $Z_i = \log[f(X_i, \theta_1)/f(X_i; \theta_0)]$ , so that the log likelihood ratio at the  $n$ th stage of sampling is  $\sum_{i=1}^n Z_i$ . The procedure at the  $n$ th stage of sampling is: accept  $H_0$  if  $\sum_{i=1}^n Z_i \leq \log B$ , accept  $H_1$  if  $\sum_{i=1}^n Z_i \geq \log A$ , continue sampling if  $\log B < \sum_{i=1}^n Z_i < \log A$ , where  $0 < B < 1 < A < \infty$  are such that

$$P_{\theta_0}[\text{Accept } H_1] = \alpha \quad \text{and} \quad P_{\theta_1}[\text{Accept } H_0] = \beta.$$

Note that the actual sample size needed for the SPRT to stop and reach a terminal decision to accept  $H_0$  or  $H_1$  is a random variable. We have to think of this random variable in the context of the sample space of infinite sequences

$$\omega = (z_1, z_2, \dots), \quad z_i = \log[f(x_i, \theta_1)/f(x_i; \theta_0)].$$

For each  $\omega$  the stopping sample size of SPRT is:  $N(\omega) = n$  if  $\log B < \sum_{i=1}^j z_i < \log A$  for  $j \leq n - 1$  and  $\sum_{i=1}^n z_i$  is either  $\leq \log B$  or  $\geq \log A$ .

To facilitate discussion, let us also define the following events in the sample space of  $\omega$ :

$$S_{n0}^* = \left\{ (z_1, z_2, \dots) : \sum_{i=1}^j z_i \in (\log B, \log A), j \leq n - 1 \text{ and } \sum_{i=1}^n z_i \leq \log B \right\},$$

and

$$S_{n1}^* = \left\{ (z_1, z_2, \dots) : \sum_{i=1}^j z_i \in (\log B, \log A), j \leq n - 1 \text{ and } \sum_{i=1}^n z_i \geq \log A \right\}$$

for  $n = 1, 2, \dots$ . These events have the following properties:

- (i)  $S_{n0}^*, S_{n1}^*, n = 1, 2, \dots$  are disjoint.
- (ii) For each  $n$ , the events  $S_{n0}^*$  and  $S_{n1}^*$  are determined by  $(Z_1, \dots, Z_n)$  only.
- (iii)  $SPRT(A, B)$  stops with exactly  $n$  observations with acceptance of  $H_0$  if  $\omega \in S_{n0}^*$  or with acceptance of  $H_1$  if  $\omega \in S_{n1}^*$ .

We now ask the following questions:

1. Does the SPRT stop with probability 1, ie, is  $P_\theta[N < \infty] = 1$ ?
2. How are the error probabilities  $\alpha, \beta$  related to the boundaries  $A, B$ ?
3. The function  $L(\theta) = P_\theta[\text{Accept } H_0] = \sum_{1 \leq n < \infty} P_\theta(S_{n0}^*)$ , called the *operating characteristic* (OC) function, which is simply 1 – power function. Obviously,  $L(\theta_0) = 1 - \alpha$  and  $L(\theta_1) = \beta$ . How to calculate  $L(\theta)$  for other values of  $\theta$ ?
4. The stopping sample size  $N$  is a random variable. It is important to know its average value  $E_\theta(N)$  for a given  $\theta$ . This average value, as a function of  $\theta$ , is called the *average sampling number* (ASN) function. How to evaluate the ASN function for a given  $\theta$ ?
5. Is the SPRT optimal in any sense?

We now deal with above issues.

### 6.11.1 SPRT Stops With Probability 1

**Theorem 6.11.1.** If  $P_\theta[Z_i = 0] < 1$ , then there exist  $c > 0$  and  $0 < r < 1$  such that  $P_\theta[N \geq n] \leq cr^n$ . In such a case,  $P_\theta[N < \infty] = 1$ .

*Remark 6.11.1.* Note that

$$P_\theta[Z_i = 0] = 1 \iff P_\theta[f(X, \theta_0) = f(X, \theta_1)] = 1.$$

If  $\theta_0$  and  $\theta_1$  are distinguishable, then  $P_{\theta_0}[f(X, \theta_0) = f(X, \theta_1)] < 1$  and  $P_{\theta_1}[f(X, \theta_0) = f(X, \theta_1)] < 1$ .

*Proof.* Since  $P_\theta[Z_i = 0] < 1$ , there exist  $\varepsilon > 0$  and  $\delta > 0$  such that either  $P_\theta[Z_i > \varepsilon] = \delta$  or  $P_\theta[Z_i < -\varepsilon] = \delta$ . We consider the former case (the other case is treated similarly).

Let  $k$  be an integer such that  $k\varepsilon > \log A - \log B$ . We now show that the desired inequality holds with  $c = (1 - \delta^k)^{-1}$  and  $r = (1 - \delta^k)^{1/k}$ . The idea of the proof is given below.

$N \geq n \iff \sum_{i=1}^l Z_i \in (\log B, \log A)$  for all  $l \leq n - 1 \implies |\sum_{i=1}^k Z_i|, \dots, |\sum_{i=(j-1)k+1}^{jk} Z_i|$  are all  $\leq \log A - \log B < k\varepsilon$ , where  $jk$  is the largest multiple of  $k$  which is  $\leq n - 1$ . Thus for  $jk + 1 \leq n \leq (j+1)k$ , using the iid property of the  $Z_i$ 's, we have

$$\begin{aligned} P_\theta[N \geq n] &\leq \left\{ P_\theta \left[ \left| \sum_{i=1}^k Z_i \right| < k\varepsilon \right] \right\}^j \leq \left\{ P_\theta \left[ \sum_{i=1}^k Z_i < k\varepsilon \right] \right\}^j \\ &\leq (1 - \delta^k)^j = (1 - \delta^k)^{-1} \{(1 - \delta^k)^{1/k}\}^{(j+1)k} \leq cr^n, \end{aligned}$$

because  $\sum_{i=1}^k Z_i < k\varepsilon \implies \{Z_1 > \varepsilon, \dots, Z_k > \varepsilon\}^c$ , so that

$$P_\theta \left[ \sum_{i=1}^k Z_i < k\varepsilon \right] \leq 1 - P_\theta[Z_1 > \varepsilon, \dots, Z_k > \varepsilon] = 1 - \delta^k.$$

Finally,

$$\begin{aligned} P_\theta[N < \infty] &= \sum_{k=1}^{\infty} P_\theta[N = k] = \lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} P_\theta[N = k] \\ &= \lim_{n \rightarrow \infty} [1 - P_\theta(N \geq n)] = \lim_{n \rightarrow \infty} (1 - cr^n) = 1 \end{aligned}$$

since  $0 < r < 1$ .

□

### 6.11.2 Error Probabilities of SPRT( $A, B$ ): Relation Between ( $A, B$ ) and $(\alpha, \beta)$

We will use [Theorem 6.11.1](#) to get

$$\begin{aligned} \alpha &= P_{\theta_0}[\text{Accept } H_1] \\ &= \sum_{n=1}^{\infty} P_{\theta_0}[\text{Stop with } n \text{ observations and accept } H_1] \\ &= \sum_{n=1}^{\infty} P_{\theta_0}[S_{n1}^*] = \sum_{n=1}^{\infty} \int_{S_{n1}^*} \prod_{i=1}^n f(x_i, \theta_0) dx_1 \cdots dx_n \\ &= \sum_{n=1}^{\infty} \int_{S_{n1}^*} \left\{ \prod_{i=1}^n f(x_i, \theta_1)/f(x_i, \theta_0) \right\}^{-1} \prod_{i=1}^n f(x_i, \theta_1) dx_1 \cdots dx_n \\ &\leq (1/A) \sum_{n=1}^{\infty} \int_{S_{n1}^*} \prod_{i=1}^n f(x_i, \theta_1) dx_1 \cdots dx_n \\ &= (1/A) P_{\theta_1}[\text{Accept } H_1] = (1/A) \{1 - P_{\theta_1}[\text{Accept } H_0]\} = \frac{1 - \beta}{A}. \end{aligned}$$

Similarly, using the fact that  $\lambda_n = \prod_{i=1}^n f(x_i, \theta_1) / f(x_i, \theta_0) \leq B$  on  $S_{n0}^*$  for all  $n$ , we have

$$\begin{aligned} 1 - \alpha &= P_{\theta_0}[\text{Accept } H_0] = \sum_{n=1}^{\infty} P_{\theta_0}[S_{n0}^*] \\ &= \sum_{n=1}^{\infty} \int_{S_{n0}^*} \lambda_n^{-1} \prod_{i=1}^n f(x_i, \theta_1) dx_1 \cdots dx_n \\ &\geq (1/B)P_{\theta_1}[\text{Accept } H_0] = \beta/B. \end{aligned}$$

Thus the Type I error probability  $\alpha(A, B)$  and Type II error probability  $\beta(A, B)$  satisfy:

$$\begin{aligned} \alpha(A, B) &\leq [1 - \beta(A, B)]/A \text{ and } 1 - \alpha(A, B) \geq \beta(A, B)/B, \text{ ie,} \\ A\alpha(A, B) + \beta(A, B) &\leq 1 \text{ and } \alpha(A, B) + \beta(A, B)/B \leq 1. \end{aligned}$$

Suppose for given  $\alpha^*, \beta^*$ , we choose  $A' = (1 - \beta^*)/\alpha^*$  and  $B' = \beta^*/(1 - \alpha^*)$ .

If the above inequalities were equalities, then we would have  $\alpha(A', B') = \alpha^*$  and  $\beta(A', B') = \beta^*$  for the resulting  $SPRT(A', B')$ . This would be the case if  $SPRT(A', B')$  always terminated with  $\lambda_N = A$  or  $\lambda_N = B$  at the stopping time  $N$ . But in reality,  $\lambda_N$  will be either  $>A$  or  $<B$  when the likelihood ratio sequence  $\{\lambda_N\}$  first goes out of the continuation region  $(B, A)$ . Taking these inequalities as equalities amounts to an approximation in which *excess over the boundaries is neglected*. Comparing  $(\alpha, \beta) = (\alpha(A', B'), \beta(A', B'))$  with  $(\alpha^*, \beta^*)$ , we see

$$\begin{aligned} \left. \begin{aligned} (i) \quad A'\alpha(A', B') + \beta(A', B') &\leq 1 \\ (ii) \quad \alpha(A', B') + \beta(A', B')/B' &\leq 1 \end{aligned} \right\} &\iff \left. \begin{aligned} (1 - \beta^*)\alpha + \alpha^*\beta &\leq \alpha^* \\ \beta^*\alpha + (1 - \alpha^*)\beta &\leq \beta^* \end{aligned} \right\} \\ &\implies \alpha + \beta \leq \alpha^* + \beta^*. \end{aligned}$$

Hence *at most one* of the error probabilities  $\alpha(A', B')$  or  $\beta(A', B')$  may exceed its intended value  $\alpha^*$  or  $\beta^*$  with the above choice of boundaries. Moreover, by (i),  $\alpha \leq 1/A' = \alpha^*/(1 - \beta^*)$ , so  $\alpha - \alpha^* \leq \alpha^*/(1 - \beta^*) - \alpha^* = \alpha^*\beta^*/(1 - \beta^*)$ , and by (ii),  $\beta \leq B' \leq \beta^*/(1 - \alpha^*)$ , so  $\beta - \beta^* \leq \beta^*/(1 - \alpha^*) - \beta^* = \alpha^*\beta^*/(1 - \alpha^*)$ .

*Summary.* If for given  $\alpha^*, \beta^*$ , we take  $A' = (1 - \beta^*)/\alpha^*$  and  $B' = \beta^*/(1 - \alpha^*)$ , then at most one of  $\alpha(A', B')$  or  $\beta(A', B')$  may exceed  $\alpha^*$  or  $\beta^*$ . If  $\alpha(A', B') > \alpha^*$ , then  $\alpha(A', B') - \alpha^* \leq \alpha^*\beta^*/(1 - \beta^*)$  and if  $\beta(A', B') > \beta^*$ , then  $\beta(A', B') - \beta^* \leq \alpha^*\beta^*/(1 - \alpha^*)$ . For example, if  $\alpha^* = \beta^* = 0.05$ , then  $\alpha - \alpha^*$  and  $\beta - \beta^*$  are both  $\leq 0.002632$  and at most one of them is positive.

In what follows, we shall use the approximate formulas:  $A \approx (1 - \beta)/\alpha$ ,  $B \approx \beta/(1 - \alpha)$ ,  $\alpha \approx (1 - B)/(A - B)$ , and  $\beta \approx (A - 1)B/(A - B)$ .

### 6.11.3 OC Function

We first prove two lemmas.

**Lemma 6.11.1.** *Let  $M(t)$  be the mgf of an rv  $Z$  and let  $T = \{t: M(t) < \infty\}$ . Then  $T$  is an interval containing 0 and  $M^{(k)}(t_0) = E[Z^k e^{t_0 Z}]$  for all  $t_0$  lying in the interior of  $T$ .*

*Proof.* Since  $M(0) = 1$ ,  $0 \in T$  and by Hölder's inequality,  $t_1$  and  $t_2$  in  $T$  implies  $\lambda t_1 + (1 - \lambda)t_2 \in T$  for all  $0 < \lambda < 1$ . Next note that by Lebesgue Dominated Convergence,

$$\begin{aligned}\frac{d}{dt}E[e^{tZ}]|_{t=t_0} &= \lim_{\delta \rightarrow 0} (1/\delta)E\left[\int e^{(t_0+\delta)z}f_Z(z) dz - \int e^{t_0z}f_Z(z) dz\right] \\ &= \int \lim_{\delta \rightarrow 0} \frac{e^{\delta z} - 1}{\delta} e^{t_0z}f_Z(z) dz = \int z e^{t_0z}f_Z(z) dz = E[Z e^{t_0 Z}],\end{aligned}$$

the dominating function being  $g$  such that  $|(e^{\delta z} - 1)/\delta| \leq g(z)$  for all  $z$  and  $\delta \neq 0$  and  $\int g(z)e^{t_0 z}f_Z(z) dz < \infty$ . Such a function exists, because for  $|\delta| < \varepsilon$ ,

$$|(e^{\delta z} - 1)/\delta| \leq \sum_{j=1}^{\infty} |\delta|^{j-1} |z|^j / j! \leq \sum_{j=1}^{\infty} \varepsilon^{j-1} |z|^j / j! < e^{\varepsilon|z|}/\varepsilon = g(z), \text{ say}$$

and

$$E\left[\varepsilon^{-1} e^{\varepsilon|Z|} e^{t_0 Z}\right] \leq \varepsilon^{-1} \{E[e^{(t_0+\varepsilon)Z}] + E[e^{(t_0-\varepsilon)Z}]\} < \infty$$

if  $\varepsilon$  is such that  $t_0 + \varepsilon$  and  $t_0 - \varepsilon$  are both in  $T$ . The proof for the higher derivatives is by induction.  $\square$

**Lemma 6.11.2.** Suppose that  $Z$  satisfies the following conditions:

- (a)  $E[Z] \neq 0$ ,
- (b)  $M(t) = E[e^{tZ}] < \infty$  for all  $t$ ,
- (c)  $P[e^Z > 1 + \delta] > 0$  and  $P[e^Z < 1 - \delta] > 0$  for some  $\delta > 0$ .

Then there exists a unique nonzero solution of the equation  $M(t) = 1$  (ie, there is a unique  $h \neq 0$  such that  $M(h) = 1$ ).

*Proof.* By condition (b) and Lemma 6.11.1,  $M''(t) = E[Z^2 e^{tZ}] > 0$ , by condition (a). Hence  $M$  is a strictly convex function. Now use condition (c) to see that: since  $M(t) \geq (1 + \delta)^t P[e^Z > 1 + \delta]$  for  $t > 0$ ,  $\lim_{t \rightarrow \infty} M(t) = \infty$ , and since  $M(t) \geq (1 - \delta)^t P[e^Z < 1 - \delta]$  for  $t < 0$ ,  $\lim_{t \rightarrow -\infty} M(t) = \infty$ . Hence  $M(t)$  has a unique minimum at some  $t_0$ . If  $t_0 = 0$ , then  $0 = M'(0) = E[Z]$  would contradict condition (a). Thus  $M(t)$  has the following properties:

- (i)  $M(0) = 1$ ,
- (ii)  $M(t)$  is strictly convex,
- (iii)  $M(t)$  has a unique minimum at  $t_0 \neq 0$ .

From these properties it follows that there exists a unique minimum  $h \neq 0$  such that  $M(h) = 1$ .  $\square$

Now consider the OC function  $L(\theta) = P_\theta[\text{Accept } H_0]$  of SPRT( $A, B$ ) for  $H_0: \theta = \theta_0$  vs  $H_1: \theta = \theta_1$ , assuming that the distribution of  $Z = \log[f(X, \theta_1)/f(X, \theta_0)]$  satisfies the conditions of Lemma 6.11.2. Then we have the unique  $h(\theta) \neq 0$  corresponding to  $\theta$ , such that

$$\begin{aligned} 1 &= E_\theta[e^{h(\theta)Z}] = E_\theta[\{f(X, \theta_1)/f(X; \theta_0)\}^{h(\theta)}] \\ &= \int \{f(X, \theta_1)/f(X; \theta_0)\}^{h(\theta)} f(x, \theta) dx \end{aligned}$$

which makes

$$f_1^*(x) = f^*(x, \theta) = \{f(X, \theta_1)/f(X; \theta_0)\}^{h(\theta)} f(x, \theta)$$

a pdf. Now let  $f_0^*(x) = f(x, \theta)$  and consider the  $SPRT(A^{h(\theta)}, B^{h(\theta)})$  for testing  $H_0^*: X \stackrel{pdf}{\sim} f_0^*$  vs  $H_1: X \stackrel{pdf}{\sim} f_1^*$ . If  $h(\theta) > 0$ , then

$$\lambda_n^* = \prod_{i=1}^n \{f_1^*(x_i)/f_0^*(x_i)\} = \left[ \prod_{i=1}^n \{f(x_i, \theta_1)/f(x_i, \theta_0)\} \right]^{h(\theta)}$$

is  $\leq B^{h(\theta)}$  or  $\geq A^{h(\theta)}$  or inside  $(B^{h(\theta)}, A^{h(\theta)})$  according as  $\lambda_n = \prod_{i=1}^n \{f(x_i, \theta_1)/f(x_i, \theta_0)\}$  is  $\leq B$  or  $\geq A$  or inside  $(B, A)$ . Hence

$$P_{H_0^*}[SPRT(A^{h(\theta)}, B^{h(\theta)}) \text{ accepts } H_0^*] = P_\theta[SPRT(A, B) \text{ accepts } H_0] = L(\theta).$$

Thus for  $h(\theta) > 0$ , neglecting excess over boundaries,

$$\begin{aligned} L(\theta) &= 1 - \alpha(A^{h(\theta)}, B^{h(\theta)}) = 1 - (1 - B^{h(\theta)})/(A^{h(\theta)} - B^{h(\theta)}) \\ &= (A^{h(\theta)} - 1)/(A^{h(\theta)} - B^{h(\theta)}). \end{aligned}$$

The same formula also holds for  $h(\theta) < 0$ . To show this, consider the correspondence between  $SPRT(B^{h(\theta)}, A^{h(\theta)})$  for testing  $H_0^*$  vs  $H_1^*$  and  $SPRT(A, B)$  for testing  $H_0$  vs  $H_1$ .

#### 6.11.4 ASN Function

**Theorem 6.11.2** (Fundamental Identity of Sequential Analysis). *Let  $S_n = \sum_{i=1}^n Z_i$  and suppose that  $P_\theta[Z_i = 0] < 1$  and  $P_\theta[|Z_i| < \infty] = 1$ . Then  $E_\theta[e^{tS_N} M(t)^{-N}] = 1$  for all  $t$  at which  $M(t) < \infty$ , where  $N$  is the stopping sample size of  $SPRT(A, B)$ .*

*Remark 6.11.2.* For each  $n$ ,

$$E_\theta[e^{tS_n}] = \{E[e^{tZ_1}]\}^n = M(t)^n, \quad \text{so } E_\theta[e^{tS_n} M(t)^{-n}] = 1.$$

The fundamental identity asserts this fact also for the random stopping time  $N$ . We shall see that the proof depends on two properties of  $N$ , namely,  $P_\theta[N < \infty] = 1$  ([Theorem 6.11.1](#)) and the fact that the event  $\{N = n\}$  depends only on  $Z_1, \dots, Z_n$ .

*Remark 6.11.3.* Note that  $Z_i = \infty$  if  $f(X_i, \theta_1) > 0$  and  $f(X_i, \theta_0) = 0$ , and  $Z_i = -\infty$  if  $f(X_i, \theta_1) = 0$  and  $f(X_i, \theta_0) > 0$ . If such  $Z_i$  are replaced by  $\log A - \log B$  and  $-(\log A - \log B)$ , respectively, then the crossing behavior of each sample sequence remains unaltered. We can, therefore, assume  $P_\theta[|Z_i| < \infty] = 1$ , without loss of generality.

*Proof.* (This proof is due to Bahadur [25].)

Let  $T = \{t: M(t) < \infty\}$  and for each  $t \in T$ , define  $p(z|t) = e^{tz}M(t)^{-1}f_Z(z|\theta)$ , where  $f_Z(z|\theta)$  is the pdf of  $Z = \log[f(X, \theta_1)/f(X; \theta_0)]$  induced by the distribution corresponding to  $f(x, \theta)$ . Then

$$\int p(z|t) dz = M(t)^{-1} \int e^{tz} f_Z(z|\theta) dz = M(t)^{-1} E_\theta[e^{tZ}] = 1,$$

showing that  $p(z|t)$  is a pdf for each  $t$  and  $p(z|0) = f_Z(z|\theta)$ . Since  $Z$  is nondegenerate under  $f_Z(z|\theta)$  by virtue of the condition  $P_\theta[Z = 0] < 1$ , it now follows that  $Z$  is nondegenerate under  $p(z|t)$  for all  $t \in T$ . Hence by [Theorem 6.11.1](#),  $P_t[N < \infty] = 1$  for all  $t \in T$ .

We now have, using the fact that the event  $\{N = n\}$  depends only on  $(Z_1, \dots, Z_n)$ ,

$$\begin{aligned} E_\theta[e^{tS_N} M(t)^{-N}] &= \sum_{n=1}^{\infty} E\left[e^{tS_N} M(t)^{-N} I(N=n)\right] \\ &= \sum_{n=1}^{\infty} \int_{N=n} \prod_{i=1}^n [e^{tz_i} M(t)^{-1} f_Z(z_i|\theta)] dz_1 \cdots dz_n \\ &= \sum_{n=1}^{\infty} \int_{N=n} \prod_{i=1}^n p(z_i|t) dz_1 \cdots dz_n \\ &= \sum_{n=1}^{\infty} P_t[N=n] = P_t[N < \infty] = 1. \end{aligned}$$

□

*Remark 6.11.4.* It can be shown that for  $t$  lying in the interior of  $T$ , differentiation can be carried under the expectation any number of times in the fundamental identity.

**Theorem 6.11.3.** Assume  $P_\theta[Z_i = 0] < 1$ ,  $P_\theta[|Z_i| < \infty] = 1$  and that  $M(t) = E_\theta[e^{tZ}] < \infty$  in a neighborhood of zero. Then

- (a)  $E_\theta[S_N] = \mu(\theta)E_\theta(N)$ ,
- (b)  $E_\theta[\{S_N - N\mu(\theta)\}^2] = \sigma^2(\theta)E_\theta(N)$ ,

where  $\mu(\theta) = E_\theta[Z]$  and  $\sigma^2(\theta) = \text{Var}_\theta(Z)$ .

*Proof.* Differentiating the fundamental identity two times under the expectation sign, we have

$$\begin{aligned} 0 &= \frac{d}{dt} E_\theta\left[e^{tS_N} M(t)^{-N}\right]_{t=0} = E_\theta\left[\left\{S_N - N\frac{M'(t)}{M(t)}\right\} e^{tS_N - N \log M(t)}\right]_{t=0} \\ &= E_\theta[S_N - NE_\theta(Z)] = E_\theta[S_N] - \mu(\theta)E_\theta[N], \text{ and} \\ 0 &= \frac{d^2}{dt^2} E_\theta\left[e^{tS_N} M(t)^{-N}\right]_{t=0} \\ &= E_\theta\left[\left\{\left(S_N - N\frac{M'(t)}{M(t)}\right)^2 - N\frac{M(t)M''(t) - M'(t)^2}{M(t)^2}\right\} e^{tS_N - N \log M(t)}\right]_{t=0} \\ &= E_\theta\left[\left[\{S_N - NE_\theta(Z)\}^2\right] - N\text{Var}_\theta(Z)\right] = E_\theta[\{S_N - N\mu(\theta)\}^2] - \sigma^2(\theta)E_\theta[N]. \end{aligned}$$

□

The ASN function is now seen to be

$$E_\theta[N] = \frac{E_\theta[S_N]}{E_\theta[Z]} \approx \frac{L(\theta) \log B + \{1 - L(\theta)\} \log A}{E_\theta[Z]}, \quad \text{if } E_\theta[Z] \neq 0,$$

where the approximation is obtained by neglecting excess over boundaries at the stopping time. Finally we state the *Optimally Property of SPRT* without proof.

Among all sequential tests with Type I and Type II error probabilities not exceeding the corresponding error probabilities of  $SPRT(A, B)$ , the ASN function of  $SPRT(A, B)$  has the smallest values at  $\theta_0$  and at  $\theta_1$ .

**Example 6.11.1.** Suppose  $X_1, X_2, \dots$  are iid  $N(\theta, \sigma^2)$ ;  $H_0: \theta = \theta_0$  and  $H_1: \theta = \theta_1 > \theta_0$ . Then

$$Z_i = \log[f(X_i, \theta_1)/f(X_i; \theta_0)] = \frac{\theta_1 - \theta_0}{\sigma^2} X_i - \frac{\theta_1^2 - \theta_0^2}{2\sigma^2}.$$

At the  $n$ th stage,

$$\begin{aligned} &\text{accept } H_0 \text{ if } \sum_{i=1}^n X_i \leq \frac{\sigma^2}{\theta_1 - \theta_0} \log B + \frac{\theta_0 + \theta_1}{2} n, \\ &\text{accept } H_1 \text{ if } \sum_{i=1}^n X_i \geq \frac{\sigma^2}{\theta_1 - \theta_0} \log A + \frac{\theta_0 + \theta_1}{2} n, \end{aligned}$$

and continue sampling if  $\sum_{i=1}^n X_i$  lies within these boundaries. Here, for  $\theta \neq (\theta_0 + \theta_1)/2$ ,

$$E_\theta[Z_i] \neq 0, \quad h(\theta) = (\theta_1 + \theta_0 - 2\theta)/(\theta_1 - \theta_0) \neq 0,$$

$$L(\theta) = \frac{A^{(\theta_1+\theta_0-2\theta)/(\theta_1-\theta_0)} - 1}{A^{(\theta_1+\theta_0-2\theta)/(\theta_1-\theta_0)} - B^{(\theta_1+\theta_0-2\theta)/(\theta_1-\theta_0)}},$$

whereas, by L'Hôpital's rule,  $L((\theta_0 + \theta_1)/2) = \log A / (\log A - \log B)$ .

**Example 6.11.2.** Suppose  $X_1, X_2, \dots$  are iid taking values 0 and 1 with probabilities  $\theta$  and  $1 - \theta$ , respectively;  $H_0: \theta = \theta_0$  and  $H_1: \theta = \theta_1 > \theta_0$ . Using the notations  $r_1 = \theta_1/\theta_0$  and  $r_2 = (1 - \theta_1)/(1 - \theta_0)$ , we have

$$Z_i = X_i \log[r_1/r_2] + \log r_2.$$

At the  $n$ th stage of sampling, accept  $H_0$  if  $\sum_{i=1}^n X_i \leq c_0 \log B + c_1 n$ , accept  $H_1$  if  $\sum_{i=1}^n X_i \geq c_0 \log A + c_1 n$  and continue sampling if  $\sum_{i=1}^n X_i$  lies within these boundaries, where  $c_0$  and  $c_1$  are easily determined. Here

$$E_\theta[Z_i] \neq 0 \text{ for } \theta \neq -\log r_2 / \log[r_1/r_2];$$

$$E_\theta[e^{hZ_i}] = \theta r_1^h + (1 - \theta) r_2^h = 1 \text{ for } \theta = \frac{1 - r_2^h}{r_1^h - r_2^h}.$$

Plot  $(\theta, L(\theta)) = \left( \frac{1 - r_2^h}{r_1^h - r_2^h}, \frac{A^h - 1}{A^h - B^h} \right)$ , using  $h$  as a parameter to obtain the OC-curve.

## 6.12 Confidence Sets

Let  $\{P_\theta: \theta \in \Omega\}$  be a family of probability distributions on  $(\mathcal{X}, \mathcal{A})$  and let  $X$  be an observable rv whose distribution belongs to this family. Let  $C$  be a mapping from  $\mathcal{X}$  into the class of all subsets of  $\Omega$ , that is, for each  $x \in \mathcal{X}$ ,  $C(x) \subset \Omega$ . We call  $C$  measurable if for each  $\theta \in \Omega$ ,  $\{x: \theta \in C(x)\} \in \mathcal{A}$ . [If  $C$  is measurable, then “ $C(X)$  covers  $\theta$ ” is an event for each  $\theta \in \Omega$ , so  $P_\theta[C(X) \text{ covers } \theta] = P_\theta[\theta \in C(X)]$  is defined.]

A measurable  $C$  is said to be a *confidence set with confidence coefficient*  $1 - \alpha$  if

$$P_\theta[\theta \in C(X)] = P_\theta[\{x: \theta \in C(x)\}] \geq 1 - \alpha \quad \text{for all } \theta \in \Omega,$$

that is,  $C(X)$  covers the true value of  $\theta$  with a probability of at least  $1 - \alpha$ .

In particular, if  $\theta$  is real and if  $C$  is such that  $C(x) = [\underline{\theta}(x), \infty)$  or  $C(x) = (-\infty, \bar{\theta}(x)]$  for each  $x$ , then  $\underline{\theta}(x)$  and  $\bar{\theta}(x)$  are called lower and upper confidence bounds, respectively, for  $\theta$  with confidence coefficient  $1 - \alpha$  and we write:  $P_\theta[\theta \in C(X)] = P_\theta[\theta \geq \underline{\theta}(X)] \geq 1 - \alpha$  for all  $\theta \in \Omega$ , or  $P_\theta[\theta \in C(X)] = P_\theta[\theta \leq \bar{\theta}(X)] \geq 1 - \alpha$  for all  $\theta \in \Omega$ .

A lower confidence bound  $\underline{\theta}(X)$  is said to be uniformly most accurate (UMA) with confidence coefficient  $1 - \alpha$  if

- (i)  $P_\theta[\theta \geq \underline{\theta}(X)] \geq 1 - \alpha$  for all  $\theta \in \Omega$ , and
- (ii)  $P_\theta[\theta \geq T(X)] \geq 1 - \alpha$  for all  $\theta \in \Omega$  implies  $P_\theta[\theta' \geq \underline{\theta}(X)] \leq P_\theta[\theta' \geq T(X)]$  for all  $\theta' < \theta$ ,

that is, among all lower confidence bounds with confidence coefficient  $1 - \alpha$ ,  $[\underline{\theta}(X), \infty)$  includes  $\theta' < \theta$  with smallest probability.

UMA upper confidence bounds analogously includes  $\theta' > \theta$  with smallest probability among all upper confidence bounds with the same confidence coefficient.

A systematic theory of confidence sets was introduced by Neyman [26]. We now discuss two methods for constructing confidence sets. Of these, the first is based on the concept of pivotal functions and the second uses a duality between acceptance regions of level  $\alpha$  tests and confidence sets with confidence coefficient  $1 - \alpha$ .

### 6.12.1 Methods Based on Pivotal Functions

**Definition 6.12.1.** A known function  $T: \mathcal{X} \times \Omega \rightarrow R$  is a pivot if the distribution of  $T(X, \theta)$  does not depend on  $\theta$  (ie, for every  $a \in R$ ,  $P_\theta[T(X, \theta) \leq a]$  is the same for all  $\theta \in \Omega$ ).

If  $T(X, \theta)$  is a pivot, then its distribution is known. This allows us to find  $c_1 < c_2$  for a given  $\alpha \in (0, 1)$  so that  $P_\theta[c_1 \leq T(X, \theta) \leq c_2] \geq 1 - \alpha$  for all  $\theta$ .

Now define  $C(x) = \{\theta \in \Omega: c_1 \leq T(x, \theta) \leq c_2\}$ . Then  $C$  is a confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$ , because  $\theta \in C(x) \iff c_1 \leq T(x, \theta) \leq c_2$ , and therefore,

$$P_\theta[\theta \in C(X)] = P_\theta[c_1 \leq T(X, \theta) \leq c_2] \geq 1 - \alpha \quad \text{for all } \theta.$$

For real-valued  $\theta$ , if  $T(x, \theta)$  is monotone in  $\theta$  for each fixed  $x$ , then  $C(x) = \{\theta: c_1 \leq T(x, \theta) \leq c_2\}$  is an interval of the form  $[\underline{\theta}(X), \bar{\theta}(X)]$ .

*Remark 6.12.1.* So far we have considered single-parameter families of probabilities. If the probabilities depend on other nuisance parameters  $\tau$  in addition to  $\theta$ , then  $C$  is a

confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$  if  $P_{\theta,\tau}[\theta \in C(X)] \geq 1 - \alpha$  for all  $(\theta, \tau) \in \Omega$ .

In such cases, a pivot  $T(x, \theta)$  should depend on  $\theta$ , but not on  $\tau$  and its distribution must be the same for all  $(\theta, \tau) \in \Omega$ . Then we can find  $c_1 < c_2$  so that  $P_{\theta,\tau}[c_1 \leq T(X, \theta) \leq c_2] \geq 1 - \alpha$  for all  $(\theta, \tau) \in \Omega$ , from which  $C(x)$  constructed as above is a confidence set with confidence coefficient  $1 - \alpha$ .

*Remark 6.12.2.* If  $T_1, \dots, T_k$  are pivots which are independent under each  $P_\theta$  in the family  $\{P_\theta : \theta \in \Omega\}$ , then any function  $g(T_1, \dots, T_k)$  is also a pivot.

**Example 6.12.1** (Location-Scale Family). Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $\sigma^{-1}f(\sigma^{-1}(x - \mu))$  where  $\mu \in R$ ,  $\sigma > 0$  and  $f$  is a known pdf.

- (a) If  $\sigma$  is known and  $\theta = \mu$  is the parameter of interest, then  $X_i - \mu, i = 1, \dots, n$  are pivots. So  $\bar{X} - \mu$  is also a pivot, from which a confidence interval for  $\mu$  is obtained as:

$$\begin{aligned} C(\mathbf{X}) &= \{\mu : c_1 \leq \bar{X} - \mu \leq c_2\} = \{\mu : \bar{X} - c_2 \leq \mu \leq \bar{X} - c_1\} \\ &= [\bar{X} - c_2, \bar{X} - c_1]. \end{aligned}$$

- (b) If  $\mu$  is known and  $\theta = \sigma$  is the parameter of interest, then  $(X_i - \mu)/\sigma, i = 1, \dots, n$  are pivots and many other pivots can be constructed as functions of these. In particular,  $S/\sigma$  is a pivot where  $S^2 = (n - 1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$  is the sample variance. Now a confidence interval for  $\sigma$  is obtained as:

$$C(\mathbf{X}) = \{\sigma : c_1 \leq S/\sigma \leq c_2\} = \{\sigma : S/c_2 \leq \sigma \leq S/c_1\} = [S/c_2, S/c_1].$$

Here  $S/\sigma$  is still a pivot if  $\mu$  is unknown and the above confidence interval for  $\sigma$  is still valid.

- (c) If  $\theta = \mu$  is the parameter of interest and  $\sigma$  is also unknown, then  $t(\mathbf{X}) = \sqrt{n}(\bar{X} - \mu)/S$ , which is called the Studentized Version of  $\sqrt{n}(\bar{X} - \mu)/\sigma$ , is a pivot, from which a confidence interval for  $\mu$  is obtained as

$$\begin{aligned} C(\mathbf{X}) &= \{\mu : c_1 \leq \sqrt{n}(\bar{X} - \mu)/S \leq c_2\} \\ &= \{\mu : \bar{X} - c_2 S/\sqrt{n} \leq \mu \leq \bar{X} - c_1 S/\sqrt{n}\} \\ &= [\bar{X} - c_2 S/\sqrt{n}, \bar{X} - c_1 S/\sqrt{n}]. \end{aligned}$$

When  $f$  is symmetric about 0, one may take  $c_1 = -c_2 = -c$  and  $C(X)$  is then of the form  $[\bar{X} - cS/\sqrt{n}, \bar{X} + cS/\sqrt{n}]$ .

**Example 6.12.2** (Ratio of Means in Bivariate Normal (Feller's Theorem)). Let  $\mathbf{X} = (X_{11}, X_{12}), \dots, (X_{n1}, X_{n2})$  be a random sample  $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu}^\top = (\mu_1, \mu_2)$  and  $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$ . The parameter of interest is  $\theta = \mu_2/\mu_1$ ,  $\mu_1 \neq 0$ . Let  $Y_i(\theta) = X_{i2} - \theta X_{i1}$ . Then  $Y_i(\theta), i = 1, \dots, n$ , are iid  $N(0, \sigma_2^2 - 2\theta\sigma_{12} + \theta^2\sigma_1^2)$ . The “sample mean” and “sample variance” of  $Y_1(\theta), \dots, Y_n(\theta)$  are  $\bar{Y}(\theta) = \bar{X}_2 - \theta\bar{X}_1$  and

$$S^2(\theta) = (n - 1)^{-1} \sum_{i=1}^n (Y_i(\theta) - \bar{Y}(\theta))^2 = S_2^2 - 2\theta S_{12} + \theta^2 S_1^2,$$

where  $S_1^2$ ,  $S_2^2$ , and  $S_{12}$  are the sample estimates of  $\sigma_1^2$ ,  $\sigma_2^2$ , and  $\sigma_{12}$ , respectively. Hence  $\sqrt{n}\bar{Y}(\theta)/S(\theta) \sim t_{n-1}$  and is therefore a pivot and a confidence set for  $\theta$  is  $C(\mathbf{X}) = \{\theta: n\bar{Y}(\theta)^2/S^2(\theta) \leq t_{n-1,\alpha/2}^2\}$ . The actual determination of  $C(\mathbf{X})$  depends on a quadratic equation in  $\theta$ .

### 6.12.2 Inverting Acceptance Regions of Tests

This method is based on a duality between hypothesis testing and confidence sets.

**Theorem 6.12.1.** *For each  $\theta'$ , let  $H_0(\theta'): \theta = \theta' \in \Omega$  and  $H_1(\theta'): \theta \in \Omega_1(\theta')$  where  $\Omega_1(\theta') \subset \Omega - \{\theta'\}$ . Then the following hold:*

- (a) *If for each  $\theta' \in \Omega$ ,  $A(\theta')$  is the acceptance region of a test for  $H_0(\theta')$  at level  $\alpha$ , then  $C(x) = \{\theta: x \in A(\theta)\}$  is a confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$ .*
- (b) *If for each  $\theta' \in \Omega$ ,  $A(\theta')$  is the acceptance region of a UMP test for  $H_0(\theta')$  vs  $H_1(\theta')$  at level  $\alpha$ , then among all confidence sets for  $\theta$  with confidence coefficient  $1 - \alpha$ ,  $C(x)$  minimizes  $P_\theta[\theta' \in C(X)]$  for all  $\theta \in \Omega_1(\theta')$ .*

*Proof.* Part (a) follows easily since

$$P_\theta[\theta \in C(X)] = P_\theta[X \in A(\theta)] \geq 1 - \alpha \quad \text{for all } \theta \in \Omega.$$

Now suppose that  $C_1(\cdot)$  is also confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$ , and let  $A_1(\theta) = \{x: \theta \in C_1(x)\}$ . Then for each  $\theta \in \Omega$ ,  $A_1(\theta)$  is the acceptance region of a level  $\alpha$  test for  $H_0(\theta)$ , because

$$P_\theta[X \in A_1(\theta)] = P_\theta[\theta \in C_1(X)] \geq 1 - \alpha \quad \text{for all } \theta.$$

Now,

$$\begin{aligned} P_\theta[\theta' \in C(X)] &= P_\theta[X \in A(\theta')] = \text{Type II error of } A(\theta') \text{ at } \theta, \\ P_\theta[\theta' \in C_1(X)] &= P_\theta[X \in A_1(\theta')] = \text{Type II error of } A_1(\theta') \text{ at } \theta. \end{aligned}$$

Since  $A(\theta')$  is the acceptance region of a UMP test for  $H_0(\theta'): \theta = \theta' \in \Omega$  vs  $H_1(\theta'): \theta \in \Omega_1(\theta')$ , it follows that for all  $\theta \in \Omega_1(\theta')$ ,

$$P_\theta[\theta' \in C(X)] \leq P_\theta[\theta' \in C_1(X)],$$

and this proves part (b). □

**Corollary.** *If  $\{P_\theta\}$  is an MLR family in  $T(x)$  and if  $F_\theta$  is the continuous cdf of  $T$  under  $P_\theta$ , then the UMA lower confidence bound for  $\theta$  with confidence coefficient  $1 - \alpha$  is given by  $\underline{\theta}(x) = \theta^*(x)$  where  $\theta^*(x)$  is the unique solution (in  $\theta$ ) of the equation  $F_\theta(T(x)) = 1 - \alpha$ .*

*Proof.* For each  $\theta'$ , the UMP level  $\alpha$  test for  $H_0(\theta'): \theta = \theta' \in \Omega$  vs  $H_1(\theta'): \theta > \theta'$  has acceptance region  $A(\theta') = \{x: T(x) \leq k(\theta')\}$  by the MLR property, where

$$P_{\theta'}[T(X) \leq k(\theta')] = F_{\theta'}(k(\theta')) = 1 - \alpha.$$

This defines a function  $k(\cdot)$  by the equation  $F_\theta(k(\theta)) = 1 - \alpha$ . The MLR property also implies that for each  $k$ ,  $F_\theta(k)$  is a strictly decreasing function of  $\theta$ , because  $F_\theta(k)$  is the Type II

error of an acceptance region  $\{T \leq k\}$  for a left-sided null hypothesis against a right-sided alternative. Hence for each  $k$ , the set  $\{\theta: F_\theta(k) \leq 1 - \alpha\}$  is an interval  $[\underline{\theta}^*(k), \infty)$ , where  $\underline{\theta}^*(k)$  is the unique solution of the equation  $F_\theta(k) = 1 - \alpha$ . Now let

$$\begin{aligned} C(x) &= \{\theta: x \in A(\theta)\} = \{\theta: T(x) \leq k(\theta)\} \\ &= \{\theta: F_\theta(T(x)) \leq F_\theta(k(\theta))\} = \{\theta: F_\theta(T(x)) \leq 1 - \alpha\}. \end{aligned}$$

Then  $C(x) = [\underline{\theta}(x), \infty)$ , where  $\underline{\theta}(x)$  is the unique solution of the equation  $F_\theta(T(x)) = 1 - \alpha$ . By part (b) of the above theorem, it now follows that  $\underline{\theta}(x)$  is UMA lower confidence bound for  $\theta$  with confidence coefficient  $1 - \alpha$ .  $\square$

**Example 6.12.3.** Let  $X = (X_1, \dots, X_n)$  be random sample from  $\text{Unif}(0, \theta)$ . We want to construct a UMA confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$ .

*Solution.* We know that the acceptance region of the UMP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1(\theta): \theta \neq \theta_0$  is

$$A(\theta_0) = \{x: \theta_0 \alpha^{1/n} \leq x_{n:n} \leq \theta_0\},$$

where  $x_{n:n} = \max(x_1, \dots, x_n)$ .

Now

$$x \in A(\theta) \iff \theta \alpha^{1/n} \leq x_{n:n} \leq \theta \iff x_{n:n} \leq \theta \leq \alpha^{-1/n} x_{n:n}.$$

It therefore follows from the above theorem that  $C(x) = [x_{n:n}, \alpha^{-1/n} x_{n:n}]$  is the UMA confidence interval for  $\theta$  with confidence coefficient  $1 - \alpha$ .

We now consider the construction of confidence sets for one parameter when there are other nuisance parameters.

Suppose that the distribution of  $X$  belongs to the family  $\{P_{\theta, \tau}\}$  where  $\theta$  is real. A confidence set for  $\theta$  with confidence coefficient  $1 - \alpha$  must satisfy:  $P_{\theta, \tau}[\theta \in C(X)] \geq 1 - \alpha$  for all  $\theta$  and  $\tau$ .

**Definition 6.12.2.**

- (i) An unbiased lower (or upper) confidence bound  $\underline{\theta}(x)$  (or  $\bar{\theta}(x)$ ) for  $\theta$  with confidence coefficient  $1 - \alpha$  must satisfy:
  - $P_{\theta, \tau}[\theta \geq \underline{\theta}(X)] \geq 1 - \alpha$  for all  $\theta, \tau$  and
  - $P_{\theta, \tau}[\theta' \geq \underline{\theta}(X)] \leq 1 - \alpha$  for all  $\theta' < \theta$  and  $\tau$ ,
  - $P_{\theta, \tau}[\theta \leq \bar{\theta}(X)] \geq 1 - \alpha$  for all  $\theta, \tau$  and
  - $P_{\theta, \tau}[\theta' \leq \bar{\theta}(X)] \leq 1 - \alpha$  for all  $\theta' > \theta$  and  $\tau$ .
- (ii) An unbiased confidence interval  $[\underline{\theta}(X), \bar{\theta}(X)]$  for  $\theta$  with confidence coefficient  $1 - \alpha$  must satisfy:
  - $P_{\theta, \tau}[\underline{\theta}(X) \leq \theta \leq \bar{\theta}(X)] \geq 1 - \alpha$  for all  $\theta, \tau$  and
  - $P_{\theta, \tau}[\underline{\theta}(X) \leq \theta' \leq \bar{\theta}(X)] \leq 1 - \alpha$  for all  $\theta' \neq \theta$  and  $\tau$ .
 Subject to these conditions, we minimize
  - $P_{\theta, \tau}[\theta' \geq \underline{\theta}(X)]$  for all  $\theta' < \theta$  and  $\tau$ , or
  - $P_{\theta, \tau}[\theta' \leq \bar{\theta}(X)]$  for all  $\theta' > \theta$  and  $\tau$ , or
  - $P_{\theta, \tau}[\underline{\theta}(X) \leq \theta' \leq \bar{\theta}(X)]$  for all  $\theta' \neq \theta$  and  $\tau$ ,

for a UMA unbiased lower confidence bound, or a UMA unbiased upper confidence bound, or a UMA unbiased confidence interval, respectively, for  $\theta$  with confidence coefficient  $1 - \alpha$ .

**Theorem 6.12.2.** *If for each  $\theta'$ ,  $A(\theta')$  is the acceptance region of a UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta'$  vs  $H_1: \theta \neq \theta'$ , then  $C(x) = \{\theta: x \in A(\theta)\}$  is a UMA unbiased confidence interval for  $\theta$  with confidence coefficient  $1 - \alpha$ .*

[Analogous correspondence holds between acceptance regions of UMP unbiased level  $\alpha$  tests for  $H_0: \theta \leq \theta'$  vs  $H_1: \theta > \theta'$  and UMA unbiased lower confidence bounds with confidence coefficient  $1 - \alpha$ , and similarly for UMA unbiased upper confidence bounds.]

**Example 6.12.4.** Let  $\mathbf{X} = (X_1, \dots, X_m)$  and  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be independent random samples from  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$ , respectively, where  $\mu_1, \mu_2, \sigma^2$  are all unknown. We want to construct a UMA unbiased confidence interval for  $\theta = \mu_1 - \mu_2$  with confidence coefficient  $1 - \alpha$ .

*Solution.* Using the notations of Example 6.9.8 in Section 6.9, note that the acceptance region of the UMP unbiased level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$  is

$$A(\theta_0) = \left\{ (\mathbf{x}, \mathbf{y}): \left| \frac{(\bar{x} - \bar{y}) - \theta_0}{SE} \right| \leq t_{m+n-2, \alpha/2} \right\},$$

where  $SE = SE(\bar{X} - \bar{Y}) = s\sqrt{1/m + 1/n}$ . Since  $(\mathbf{x}, \mathbf{y}) \in A(\theta_0)$  is equivalent to

$$(\bar{x} - \bar{y}) - t_{m+n-2, \alpha/2}SE \leq \theta_0 \leq (\bar{x} - \bar{y}) + t_{m+n-2, \alpha/2}SE,$$

the UMA unbiased confidence interval for  $\theta$  with confidence coefficient  $1 - \alpha$  is

$$C(\mathbf{x}, \mathbf{y}) = [(\bar{x} - \bar{y}) - t_{m+n-2, \alpha/2}SE, (\bar{x} - \bar{y}) + t_{m+n-2, \alpha/2}SE].$$

## Exercises

**6.1.** Let  $f(\mathbf{x}, \theta) = \prod_{i=1}^n [c(\theta)h(x_i)I_{(-\infty, \theta)}(x_i)]$  and define  $x_{n:n} = \max(x_1, \dots, x_n)$ .

- (a) Show that  $\{f(\mathbf{x}, \theta): \theta \in \mathbb{R}\}$  is an MLR family.
- (b) Express the joint pdf of a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  from  $\text{Unif}(0, \theta)$  as a special case of this family.
- (c) Show that the test

$$\phi(\mathbf{x}) = \begin{cases} 1 & \text{if } x_{n:n} > \theta_0 \\ 0 & \text{if } x_{n:n} \leq \theta_0 \end{cases}$$

is a UMP level  $\alpha$  test for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$ , but that  $\phi$  is inadmissible under the 0 – 1 loss function (ie, there exists a test  $\psi$  whose risk function under the 0 – 1 loss function satisfies  $R(\theta, \psi) \leq R(\theta, \phi)$  for all  $\theta$ ), with strict inequality holding for some  $\theta$ .

- (d) Show that

$$\phi(\mathbf{x}) = \begin{cases} 1 & \text{if } x_{n:n} > \theta_0 \text{ or } x_{n:n} \leq b \\ 0 & \text{if } b < x_{n:n} \leq \theta_0, \end{cases}$$

where  $b = \theta_0\alpha^{1/n}$  is a UMP level  $\alpha$  test for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .

- 6.2.** Let  $X$  be the number of successes in  $n$  independent trials with  $P[\text{success}] = \theta$ . Let  $\phi(x)$  be the UMP level  $\alpha$  test for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$ .
- (a) For  $n = 6$ ,  $\theta_0 = 0.25$  and  $\alpha = 0.05, 0.10, 0.20$ , find  $c$  and  $\gamma$  for  $\phi$  and find the powers of these tests at  $\theta = 0.40$ .
  - (b) Use normal approximation in order to find the smallest  $n$  required for the UMP test at level  $\alpha = 0.05$  for  $H_0: \theta \leq 0.25$  vs  $H_1: \theta > 0.25$  to attain power  $\beta(0.4) = 0.90$ .
- 6.3.** Let  $X_1, \dots, X_n$  be a random sample from an exponential distribution with mean  $\theta$ .
- (a) For  $n = 5$ , find the UMP test for  $H_0: \theta \leq 1$  vs  $H_1: \theta > 1$  at level  $\alpha = 0.1$ . [The exponential rv with mean 2 is the same as a  $\chi^2$  with 2 df.]
  - (b) For the test obtained in part (a), calculate the power  $\beta(2)$  at  $\theta = 2$ .
  - (c) Find the smallest  $n$  so that the UMP test at level  $\alpha = 0.1$  has  $\beta(2) = 0.8$ .  
[For parts (b) and (c), use integration by parts to show that  $\int_c^\infty x^n \exp(-x) dx/n! = P[W \leq n]$ , where  $W \sim \text{Poi}(c)$ .]
- 6.4.** A box contains  $N$  manufactured items of which an unknown number  $\theta$  are defective and the other  $N - \theta$  are good. Let  $X$  denote the number of defective items in a random sample of  $n$  items drawn without replacement from the box. Then  $X$  has pmf
- $$f(x, \theta) = P_\theta[X = x] = \binom{\theta}{x} \binom{N - \theta}{n - x} / \binom{N}{n}, \quad x = \max(0, n + \theta - N), \dots, \min(n, \theta).$$
- (a) Show that  $\{f(x, \theta): \theta = 0, 1, \dots, N\}$  is an MLR family. [For  $\theta_1 < \theta_2$ , write  $f(x, \theta_2)/f(x, \theta_1) = \prod_{j=\theta_1}^{\theta_2-1} \{f(x, j+1)/f(x, j)\}$  and examine how  $f(x, j+1)/f(x, j)$  changes with  $x$ .]
  - (b) For a specified integer  $\theta_0$ , write down the UMP test at a given level  $\alpha$  for  $H_0: \theta \leq \theta_0$  vs  $H_1: \theta > \theta_0$ . Explain how the constants involved in the UMP test are determined.
- 6.5.** Let  $X_1, \dots, X_n$  be a random sample from  $\text{Unif}(\theta, \theta + 1)$ .
- (a) Show that  $(T_1, T_2) = (X_{n:1}, X_{n:n})$  are jointly sufficient for  $\theta$  and find the joint distribution of  $(T_1, T_2)$ .
  - (b) Show that the UMP test at level  $\alpha$  for  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  is of the form
- $$\phi(t_1, t_2) = \begin{cases} 0 & \text{if } t_1 < k \text{ and } t_2 < 1 \\ 1 & \text{otherwise} \end{cases}, \quad \text{where } k = 1 - \alpha^{1/n}.$$
- 6.6.** Let  $P_0$  and  $P_1$  be two probability distributions with pdf's  $p_0$  and  $p_1$ , respectively. Suppose that under  $P_0$ , the likelihood ratio  $T = p_1(X)/p_0(X)$  has a pdf which is everywhere positive. For  $0 < \alpha < 1$ , let  $\{x: T(x) \geq k(\alpha)\}$  denote the critical region of an N-P test of size  $\alpha$  for  $H_0: P = P_0$  vs  $H_1: P = P_1$  and let  $\beta(\alpha) = P_1[T \geq k(\alpha)]$ . Show that  $\beta'(\alpha) = k(\alpha)$ .
- 6.7.** Suppose  $P_0 \neq P_1$  are probabilities on  $(\mathcal{X}, \mathcal{A})$  and  $X_1, \dots, X_n$  are independent samples from  $(\mathcal{X}, \mathcal{A}, P)$  where  $P$  is either  $P_0$  or  $P_1$ . We want to test  $H_0: P = P_0$  vs  $H_1: P = P_1$ . Show that there exists a sequence of tests  $\{\phi_n\}$ , each based on  $(X_1, \dots, X_n)$  such that

$$\lim_{n \rightarrow \infty} E_{P_0}[\phi_n(X_1, \dots, X_n)] = 0 \text{ and } \lim_{n \rightarrow \infty} E_{P_1}[\phi_n(X_1, \dots, X_n)] = 1,$$

that is, the Type I error probability converges to zero and the power converges to 1.

[Hint: Let  $p_i$  denote the pdf/pmf corresponding to  $P_i$ ,  $i = 0, 1$ . Then for any  $\alpha$ , the MP test at level  $\alpha$  for  $H_0$  vs  $H_1$  based on  $(X_1, \dots, X_n)$  is roughly (except for discreteness) of the form:

$$\phi_n(X_1, \dots, X_n) = I_{[k_n, \infty)} \left[ n^{-1} \sum_{i=1}^n \log(p_1(X_i)/p_0(X_i)) \right],$$

where  $k_n$  is chosen according to  $\alpha$ . Let

$$m_i = E_{P_i}[\log(p_1(X_i)/p_0(X_i))], \quad i = 0, 1.$$

Use Jensen's inequality to show that  $m_0 < 0 < m_1$ , choose  $k_n = (m_0 + m_1)/2$  in the definition of  $\phi_n$ . Now use the SLLN.]

- 6.8. Let  $X \sim \text{Bin}(n, p)$  with  $n = 10$  and  $p$  unknown. Find the UMP unbiased tests at level  $\alpha = 0.10$  for
  - (a)  $H_0: p = 0.2$  vs  $H_1: p \neq 0.2$ , (b)  $H_0: p = 0.5$  vs  $H_1: p \neq 0.5$ , (c)  $H_0: 0.4 \leq p \leq 0.6$  vs  $H_1: p \notin [0.4, 0.6]$ .
- 6.9. Let  $X$  and  $Y$  be independent Poisson rv's with means  $\lambda$  and  $\mu$ , respectively. Construct UMP unbiased tests at level  $\alpha = 0.1$  over the set  $\{(x, y): x + y = 8\}$  for
  - (a)  $H_0: \lambda \leq \mu$  vs  $H_1: \lambda > \mu$ , (b)  $H_0: \lambda = \mu$  vs  $H_1: \lambda \neq \mu$ .
- 6.10. Let  $X_1, \dots, X_n$  be a random sample from an exponential distribution with mean  $\theta$ . Find the UMP unbiased level  $\alpha$  test for  $H_0: \theta = 2$  vs  $H_1: \theta \neq 2$ .
- 6.11. Let  $X$  be a random sample of size 1 from the beta distribution  $Be(\theta, 1)$  with the pdf  $f(x, \theta) = \theta x^{\theta-1}$ ,  $0 < x < 1$ . Find the UMP unbiased test at level  $\alpha = 0.1$  for  $H_0: \theta = 1$  vs  $H_1: \theta \neq 1$  and determine the critical value.
- 6.12. Use Basu's Theorem to show the following:
  - (a) If  $X_1, \dots, X_n$  is a random sample from  $N(0, \sigma^2)$ , then  $\sum X_i / \sqrt{\sum X_i^2}$  and  $\sum X_i^2$  are independent.
  - (b) If  $(X_1, Y_1), \dots, (X_n, Y_n)$  are iid  $N_2(\mu, \Sigma)$ , where  $\mu = (\mu_1, \mu_2)$  and  $\Sigma$  is a diagonal matrix with diagonal entries  $\sigma_1^2$  and  $\sigma_2^2$ , then  $\sum (X_i - \bar{X}_n)^2$ ,  $\sum (Y_i - \bar{Y}_n)^2$  and the sample correlation  $r$  are mutually independent.
- 6.13. Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a random sample from  $Unif(\theta_1, \theta_2)$ . Sufficient statistic for  $(\theta_1, \theta_2)$  in  $\mathbf{X}$  is  $(T_1, T_2) = (X_{n:1}, X_{n:n})$ .
  - (a) Show that  $T_1$  given  $T_2 = t_2$  is distributed as  $U_{n-1:1}$  where  $(U_1, \dots, U_{n-1})$  is a random sample from  $Unif(\theta_1, t_2)$ .
  - (b) Construct a UMP level  $\alpha$  test for  $H_0: \theta_1 \leq 0$  vs  $H_1: \theta_1 > 0$  conditionally, given  $T_2 = t_2$ . Call this test  $\phi_{t_2}(t_1)$ .
  - (c) Show that  $\phi(t_1, t_2) = \phi_{t_2}(t_1)$  is unconditionally a UMP level  $\alpha$  test for  $H_0$  vs  $H_1$ .
- 6.14. Let  $X$  and  $Y$  be two independent exponential random variables with means  $1/\lambda$  and  $1/\mu$ , respectively. Find UMP unbiased tests at level  $\alpha = 0.2$  for

- (a)  $H_0: \lambda \leq \mu + 1$  vs  $H_1: \lambda > \mu + 1$ , (b)  $H_0: \lambda = \mu$  vs  $H_1: \lambda \neq \mu$ , (c)  $H_0: \lambda \geq 2\mu$  vs  $H_1: \lambda < 2\mu$ .

**6.15.** Let  $X_1, X_2$  be independent rv's with pmf's

$$f_{X_i}(x, \theta_i) = \theta_i(1 - \theta_i)^{x-1}, \quad x = 1, 2, \dots, \text{ where } 0 < \theta_i < 1.$$

Find UMP unbiased tests at level  $\alpha = 0.2$  for

- (a)  $H_0: \theta_1 \leq \theta_2$  vs  $H_1: \theta_1 > \theta_2$ , (b)  $H_0: \theta_1 = \theta_2$  vs  $H_1: \theta_1 \neq \theta_2$ .

- 6.16.** Let  $X = (X_1, \dots, X_m)$  and  $Y = (Y_1, \dots, Y_n)$  be independent samples from  $N(\mu_1, \sigma_1^2 = 1)$  and  $N(\mu_2, \sigma_2^2 = 2)$ , respectively. We want to test  $H_0: \mu_1 \leq \mu_2$  vs  $H_1: \mu_1 > \mu_2$ . Derive the UMP unbiased level  $\alpha$  test for  $H_0$  vs  $H_1$  by first expressing it in a conditional form (in terms of appropriate sufficient statistics), and then unconditionally in terms of  $\bar{X}_m - \bar{Y}_n$ .
- 6.17.** Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be a random sample from a bivariate normal distribution with  $E[X_1] = E[Y_1] = \mu$ ,  $\text{Var}[X_1] = \text{Var}[Y_1] = \sigma^2$  and  $\text{Cov}[X_1, Y_1] = \sigma^2/2$ , with  $\mu$  and  $\sigma^2$  unknown. Find the UMP unbiased level  $\alpha$  test for  $H_0: \mu = 0$  vs  $H_1: \mu \neq 0$ . [Transform:  $U_i = (X_i + Y_i)/\sqrt{3}$  and  $V_i = X_i - Y_i$ .]
- 6.18.** Let  $X_1, X_2, \dots$  be sequentially observed independent samples from a normal distribution with unknown mean  $\theta$  and known SD  $\sigma = 5$ . We want to test  $H_0: \theta = 0$  vs  $H_1: \theta = 2$  at level  $\alpha = 0.01$ , holding the probability of Type II error probability at  $\beta = 0.05$ .
- (a) Find the approximate values of  $A, B$  in  $SPRT(A, B)$  needed for this purpose and describe the procedure in terms of cumulative sums  $\sum_{i=1}^n X_i$ ,  $n = 1, 2, \dots$
  - (b) Evaluate the OC function  $L(\theta) = P_\theta[\text{accept } H_0]$  for  $\theta = 3$ .
  - (c) Evaluate the ASN function  $E_\theta(N)$  for  $\theta = 0$  and  $\theta = 2$ .
  - (d) Find the smallest sample size  $n(\alpha, \beta)$  needed for a fixed sample size test with  $\alpha = 0.01$  and  $\beta \leq 0.05$  in this problem, and compare  $n(\alpha, \beta)$  with the numbers obtained in (c).
- 6.19.** Do problem 18 when the  $X_i$ 's are sequentially observed independent samples from
- (i) Exponential distribution with mean  $\theta$ ,
  - (ii) Poisson distribution with mean  $\theta$ , and we want to test  $H_0: \theta = 1$  vs  $H_1: \theta = 2$  at level  $\alpha = 0.05$ , holding  $\beta = 0.10$ .
- 6.20.** Let  $X_1, \dots, X_{10}$  be a random sample from an exponential distribution with mean  $\theta$ . Find the UMA lower confidence bound for  $\theta$  with confidence coefficient  $1 - \alpha = 0.95$  based on  $(X_1, \dots, X_{10})$ , using the table for the  $\chi^2$ -distribution.
- 6.21.** Let  $X_1, \dots, X_n$  be iid following the Weibull distribution with the pdf

$$f(x, \lambda) = \lambda c x^{c-1} \exp(-\lambda x^c), \quad x > 0,$$

where  $c$  is known but  $\lambda$  is unknown. Show that the UMA upper confidence bound for  $\theta = 1/\lambda$  with confidence coefficient  $1 - \alpha$  is given by  $\bar{\theta} = 2 \sum X_i^c / \chi_{2n}^2(\alpha)$  where  $P[\chi_{2n}^2 \leq \chi_{2n}^2(\alpha)] = \alpha$ . [Hint: Find the distribution of  $X_i^c$ .]

- 6.22.** Let  $X$  and  $Y$  be independent exponential random variables with means  $1/\lambda$  and  $1/\mu$ , respectively. Construct a UMA unbiased confidence interval with confidence

coefficient  $1 - \alpha$  for  $\theta = \lambda/\mu$  by inverting the acceptance regions of UMP unbiased level  $\alpha$  tests for  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0, \theta_0 > 0$ .

- 6.23.** Four independent experiments were carried out to test whether a new method of irrigation would provide better yield of a certain crop. The data from all four experiments were analyzed to test the same hypothesis  $H_0$  that the new method is no better than the old vs the alternative  $H_1$  that the new method is better. The results of these experiments are summarized by four  $t$ -statistics given below with their respective df's

$$t_1 = 1.50, df = 10; t_2 = 2.15, df = 16; t_3 = 1.80, df = 19; t_4 = 1.30, df = 25.$$

- (a) Calculate the  $p$ -values of the four statistics
- (b) Calculate the  $P_\lambda$ -statistic combining the results of all four experiments.
- (c) Comment.

# Methods Based on Likelihood and Their Asymptotic properties

## 7.1 Asymptotic Properties of the MLEs: Consistency and Asymptotic Normality

On a historical note, the maximum likelihood estimators (MLEs) were introduced by R.A. Fisher in early 1920s, which he claimed to be “better” than the method of moments estimators used widely, especially by K. Pearson. To justify the superiority of the MLEs, Fisher used the concepts of consistency (converging to the true parameter in probability), sufficiency (capturing everything relevant in the sample), and efficiency (attaining the smallest possible variance among all unbiased estimators, which led to the definition of Fisher-information). This led to the foundation of the theory of statistical inference in the area of estimation. Asymptotic properties of the MLEs will be discussed in this section.

Let  $X_1, \dots, X_n, \dots$  be iid with pdf/pmf  $f(x; \theta_0)$  in the family  $\{f(x; \theta) : \theta \in \Theta\}$ , where  $\theta_0$ , the unknown value of the parameter is an interior point of  $\Theta$ . All probability statements (including  $\xrightarrow{P}$ ,  $\xrightarrow{\mathcal{L}}$ ,  $o_P$ ,  $O_P$ ) and all expectations, variances, and covariances are with respect to  $f(\cdot, \theta_0)$ , unless stated otherwise. We assume throughout that the family  $\{f(x, \theta) : \theta \in \Theta\}$  satisfies the identifiability condition introduced in [Section 5.5.1](#).

The MLE of  $\theta_0$  based on  $(X_1, \dots, X_n)$  is denoted by

$$\hat{\theta}_n = \arg \max_{t \in \Theta} \prod_{i=1}^n f(X_i, t) = \arg \max_{t \in \Theta} \sum_{i=1}^n l(X_i, t), \quad (1)$$

where  $l(x, t) = \log f(x, t)$ .

Using the notations of [Section 5.2.1](#), we write  $\log f(x, t) = l(x, t)$  as in the above paragraph and also

$$\begin{aligned} \frac{\partial f(x, \theta)}{\partial \theta} &= \dot{f}(x, \theta), \quad \frac{\partial^2 f(x, \theta)}{\partial \theta^2} = \ddot{f}(x, \theta), \\ \frac{\partial l(x, \theta)}{\partial \theta} &= \dot{l}(x, \theta) \text{ and } \frac{\partial^2 l(x, \theta)}{\partial \theta^2} = \ddot{l}(x, \theta) \end{aligned} \quad (2a)$$

in the single-parameter case when  $\Theta \subset \mathbb{R}$  and

$$\begin{aligned}\frac{\partial f(x, \theta)}{\partial \theta_r} &= \dot{f}_r(x, \theta), \quad \frac{\partial^2 f(x, \theta)}{\partial \theta_r \partial \theta_s} = \ddot{f}_{rs}(x, \theta), \\ \frac{\partial l(x, \theta)}{\partial \theta_r} &= \dot{l}_r(x, \theta) \text{ and } \frac{\partial^2 l(x, \theta)}{\partial \theta_r \partial \theta_s} = \ddot{l}_{rs}(x, \theta)\end{aligned}\tag{2b}$$

in the multiparameter case when  $\Theta \subset \mathbb{R}^k$ .

The following conditions were essentially introduced by Cramér [18, p. 500–501].

Regularity conditions (Cramér): single-parameter case.

In the notations of Eqs. (2a) and (2b):

1.  $\dot{f}(x, t), \ddot{f}(x, t)$  exist for all  $(x, t)$  and there exists a nonnegative function  $g(x)$  with  $\int g < \infty$  such that  $|\dot{f}(x, t)|$  and  $|\ddot{f}(x, t)|$  are bounded above by  $g(x)$  for all  $(x, t)$ .
2. There exist nonnegative functions  $H(x, \theta_0)$  and  $\phi(\varepsilon)$  such that

$$\sup_{|t-\theta_0| \leq \varepsilon} |\ddot{l}(x, t) - \ddot{l}(x, \theta_0)| \leq H(x, \theta_0)\phi(\varepsilon),$$

where  $\lim_{\varepsilon \rightarrow 0} \phi(\varepsilon) = 0$  and  $E[H(X, \theta_0)] < \infty$ .

3.  $0 < -E[\ddot{l}(X, \theta_0)] = I(\theta_0) < \infty$ .

By dominated convergence, Condition 1 allows differentiation of  $\int f(x, \theta) dx$  twice with respect to  $\theta$  under the integral. Since  $\int f(x, \theta) dx = 1$  for all  $\theta$ , we have

$$\begin{aligned}0 &= \frac{d}{d\theta} \int f(x, \theta) dx|_{\theta=\theta_0} = \int \dot{f}(x, \theta_0) dx \\ &= \int \dot{l}(x, \theta_0) f(x, \theta_0) dx = E[\dot{l}(X, \theta_0)], \text{ and} \\ 0 &= \frac{d^2}{d\theta^2} \int f(x, \theta) dx|_{\theta=\theta_0} = \int \ddot{f}(x, \theta_0) dx \\ &= E[\ddot{l}(X, \theta_0)] + E[\dot{l}(X, \theta_0)]^2.\end{aligned}$$

Thus

$$E[\dot{l}(X, \theta_0)] = 0 \text{ and } I(\theta_0) = E[-\ddot{l}(X, \theta_0)] = \text{Var}[\dot{l}(X, \theta_0)].\tag{3a}$$

Regularity conditions in the multiparameter case: Condition 1 should hold for  $\dot{f}_r(x, t)$  and  $\ddot{f}_{rs}(x, t)$  for all  $(x, t)$  and for all  $r$  and  $s$ , Condition 2 should hold for  $\dot{l}_{rs}(x, t)$  for all  $r$  and  $s$ , and in Condition 3 we need all elements of the information matrix  $I(\theta_0) = ((I_{rs}(\theta_0)))$  to exist and  $I(\theta_0)$  to be positive definite.

By the same argument as in the single-parameter case, it follows from Condition 1 by dominated convergence, that

$$\begin{aligned}E[\dot{l}_r(X, \theta_0)] &= 0 \text{ and} \\ I_{rs}(\theta_0) &= E[-\ddot{l}_{rs}(X, \theta_0)] = \text{Cov}[\dot{l}_r(X, \theta_0), \dot{l}_s(X, \theta_0)]\end{aligned}\tag{3b}$$

for all  $r$  and  $s$ .

**Theorem 7.1.1.** *Under regularity Condition 1, the MLE  $\hat{\theta}_n$  is consistent (ie,  $\hat{\theta}_n \xrightarrow{P} \theta_0$  as  $n \rightarrow \infty$ ).*

*Proof.* As shown in [Section 5.5.1](#), the identifiability condition implies that  $E[l(X, t)]$  has a unique maximum at  $t = \theta_0$ . In the single-parameter case, we therefore have for any  $\delta > 0$

$$E[l(X, \theta_0)] > \max\{E[l(X, \theta_0 - \delta)], E[l(X, \theta_0 + \delta)]\}.$$

Hence for a given  $\delta > 0$ , there exists  $\varepsilon > 0$  such that

$$E[l(X, \theta_0)] - E[l(X, t)] > 2\varepsilon \quad \text{for } t = \theta_0 \pm \delta.$$

Since  $n^{-1} \sum l(X_i, t) \xrightarrow{P} E[l(X, t)]$  for all  $t$ , it follows that

$$\lim_{n \rightarrow \infty} P\left[n^{-1} \sum l(X_i, \theta_0) - n^{-1} \sum l(X_i, t) > \varepsilon\right] = 1 \quad \text{for } t = \theta_0 \pm \delta.$$

But  $\dot{l}(x, t)$  exist in a neighborhood of  $\theta_0$ , so with probability tending to 1 as  $n \rightarrow \infty$ , the equation

$$g_n(t) = n^{-1} \sum_{i=1}^n \dot{l}(X_i, t) = 0$$

has a solution in the intervals  $(\theta_0 - \delta, \theta_0 + \delta)$  for arbitrary  $\delta > 0$ ; that is, there is a sequence of solutions  $\{\hat{\theta}_n\}$  of the equation  $g_n(t) = 0$  needed for (1) converges to  $\theta_0$  in probability as  $n \rightarrow \infty$ . The MLE is consistent in this sense.

In the multiparameter case, we modify the above proof by replacing  $\theta_0 \pm \delta$  with  $S(\theta_0, \delta) = \{t: \|t - \theta_0\| \leq \delta\}$  and argue in the same way, assuming that the regularity conditions hold.  $\square$

**Theorem 7.1.2** (Asymptotic Normality of MLEs: Single-Parameter Case). *Under regularity Conditions 1, 2, 3, the MLE  $\hat{\theta}_n$  is asymptotically normal (ie,  $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta_0))$ ).*

*Proof.* Let  $\{\hat{\theta}_n\}$  be the sequence described in the proof of [Theorem 7.1.1](#) with  $\hat{\theta}_n = \theta_0 + o_P(1)$  and  $\sum_{i=1}^n \dot{l}(X_i, \hat{\theta}_n) = 0$ . Then

$$0 = n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \hat{\theta}_n) = n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta_0) + \sqrt{n}(\hat{\theta}_n - \theta_0)n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_n),$$

where  $\tilde{\theta}_n = \theta_0 + \lambda(\hat{\theta}_n - \theta_0)$ ,  $0 \leq \lambda \leq 1$ . Hence

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &= \frac{n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta_0)}{-n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_n)} \\ &= \frac{n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta_0)}{-n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \theta_0) + n^{-1} \sum_{i=1}^n R_n(X_i)}. \end{aligned} \tag{4}$$

On the right-hand side of this expression, the numerator  $n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta_0) \xrightarrow{\mathcal{L}} N(0, I(\theta_0))$  since  $\dot{l}(X_i, \theta_0)$ ,  $i = 1, 2, \dots$  are iid with mean 0 and variance  $I(\theta_0)$  by Eq. (3a), while in the denominator,

$$-n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \theta_0) \xrightarrow{P} E[-\ddot{l}(X, \theta_0)] = I(\theta_0).$$

Finally,

$$\begin{aligned} |R_n(X_i)| &= |\ddot{l}(X_i, \theta_0 + \lambda(\hat{\theta}_n - \theta_0)) - \ddot{l}(X_i, \theta_0)| \\ &\leq \sup_{|\theta' - \theta_0| \leq |\hat{\theta}_n - \theta_0|} |\ddot{l}(X_i, \theta') - \ddot{l}(X_i, \theta_0)| \leq \varphi(|\hat{\theta}_n - \theta_0|) H(X_i, \theta_0), \end{aligned}$$

by Condition 3, so

$$\begin{aligned} \left| n^{-1} \sum_{i=1}^n R_n(X_i) \right| &\leq \varphi(|\hat{\theta}_n - \theta_0|) n^{-1} \sum_{i=1}^n H(X_i, \theta_0) \\ &= o_P(1)\{O(1) + o_P(1)\} = o_P(1), \end{aligned}$$

since  $|\hat{\theta}_n - \theta_0| = o_P(1)$  implies  $\varphi(|\hat{\theta}_n - \theta_0|) = o_P(1)$  and  $n^{-1} \sum_{i=1}^n H(X_i, \theta_0) = E[H(X, \theta_0)] + o_P(1) = O(1) + o_P(1)$ . Putting all this together in Eq. (4), we see that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta_0)}{I(\theta_0) + o_P(1) + o_P(1)} \xrightarrow{\mathcal{L}} N(0, 1/I(\theta_0)).$$

□

Suppose a sequence of unbiased estimators  $T_n^* = T_n^*(X_1, \dots, X_n)$  of  $\theta_0$  has  $\text{Var}[T_n^*] = 1/\{nI(\theta_0)\}$  which is the Cramér-Rao lower bound for unbiased estimators (under regularity conditions). If for a sequence of unbiased estimators, we define

$$e_n(T_n) = \frac{\text{C-R Lower Bound}}{\text{Var}[T_n]} = \frac{1/\{nI(\theta_0)\}}{\text{Var}[T_n]}$$

as the efficiency of  $T_n$ , then  $T_n^*$  described above has efficiency 1.

The large sample analog of this is to make comparison among all *consistent estimators* in terms of their *asymptotic variance*. Since the MLE  $\hat{\theta}_n$  has variance  $1/\{nI(\theta_0)\}$  in an asymptotic sense, we can say that the asymptotic efficiency of the MLE is 1, or simply that the MLE is asymptotically efficient. However, the justification of the last statement needs much deeper analysis as will be seen later.

**Definition 7.1.1.** Any estimator  $T_n$  with  $\sqrt{n}(T_n - \theta_0) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta_0))$  is said to be a *best asymptotically normal* (BAN) estimator. (This is because of asymptotic normality and the property of asymptotic efficiency.)

## Finding a BAN Estimator by the Newton-Raphson Method

In some situations, the likelihood equation does not have a closed form solution, so we cannot obtain the MLE  $\hat{\theta}_n$  explicitly. However, in such a case, we can often find an estimator which is consistent. Starting with such an estimator  $\tilde{\theta}_{n0} = \theta_0 + o_P(1)$  as an initial estimator we can use the Newton-Raphson method. After one iteration, we have

$$\tilde{\theta}_{n1} = \tilde{\theta}_{n0} + \frac{n^{-1} \sum_{i=1}^n \dot{l}(X_i, \tilde{\theta}_{n0})}{-n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_{n0})}.$$

By argument given in the course of the proof of [Theorem 7.1.2](#), it can be shown that if  $T_n = \theta_0 + o_P(1)$ , then  $\left| n^{-1} \sum_{i=1}^n \ddot{l}(X_i, T_n) - n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \theta_0) \right| = o_P(1)$ . Since  $\tilde{\theta}_{n0} = \theta_0 + o_P(1)$ , we then have

$$\begin{aligned} 0 &= n^{-1} \sum_{i=1}^n \dot{l}(X_i, \hat{\theta}_n) = n^{-1} \sum_{i=1}^n \dot{l}(X_i, \tilde{\theta}_{n0}) + (\hat{\theta}_n - \tilde{\theta}_{n0}) n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_{n0}) \\ &\quad + o_P(1)(\hat{\theta}_n - \tilde{\theta}_{n0}), \end{aligned}$$

so that

$$\hat{\theta}_n = \tilde{\theta}_{n0} + \frac{n^{-1} \sum_{i=1}^n \dot{l}(X_i, \tilde{\theta}_{n0})}{-n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \tilde{\theta}_{n0})} + o_P(1)(\hat{\theta}_n - \tilde{\theta}_{n0}) = \tilde{\theta}_{n1} + o_P(1)(\hat{\theta}_n - \tilde{\theta}_{n0}).$$

Hence

$$\sqrt{n}(\tilde{\theta}_{n1} - \theta_0) = \sqrt{n}(\hat{\theta}_n - \theta_0) + o_P(1)\sqrt{n}(\hat{\theta}_n - \tilde{\theta}_{n0}) = \sqrt{n}(\hat{\theta}_n - \theta_0) + o_P(1)$$

under additional condition that  $\sqrt{n}(\hat{\theta}_n - \tilde{\theta}_{n0}) = O_P(1)$  (ie,  $\tilde{\theta}_{n0} - \theta_0 = O_P(n^{-1/2})$ ). In such a case,  $\sqrt{n}(\tilde{\theta}_{n1} - \theta_0)$  has the same asymptotic distribution as that of  $\sqrt{n}(\hat{\theta}_n - \theta_0)$ , so  $\tilde{\theta}_{n1}$  is a BAN estimator.

In the multiparameter case,  $\boldsymbol{\theta}_0$  and  $\mathbf{t}$  are  $k$ -dim column vectors with  $\theta_{0r}$  and  $t_r$  as their  $r$ th coordinates, and

$$\dot{\mathbf{l}}(x, \mathbf{t})^T = (\dot{l}_1(x, \mathbf{t}), \dots, \dot{l}_k(x, \mathbf{t})), \quad \ddot{\mathbf{l}}(x, \mathbf{t}) = ((\ddot{l}_{rs}(x, \mathbf{t}))).$$

By regularity conditions,  $E[\dot{l}_r(X, \boldsymbol{\theta}_0)] = 0$  and

$$E[-\ddot{l}_{rs}(X, \boldsymbol{\theta}_0)] = \text{Cov}[\dot{l}_r(X, \boldsymbol{\theta}_0), \dot{l}_s(X, \boldsymbol{\theta}_0)] = I_{rs}(\boldsymbol{\theta}_0)$$

as shown in Eq. (3b). In matrix notation,

$$E[\dot{\mathbf{l}}(X, \boldsymbol{\theta}_0)] = \mathbf{0} \text{ and } E[-(\ddot{l}_{rs}(X, \boldsymbol{\theta}_0))] = ((I_{rs}(\boldsymbol{\theta}_0))) = \mathbf{I}(\boldsymbol{\theta}_0).$$

Let

$$A_{nr} = n^{-1/2} \sum_{i=1}^n \dot{l}_r(X_i, \boldsymbol{\theta}_0) \text{ and } \mathbf{A}_n^T = (A_{n1}, \dots, A_{nk}).$$

Then  $\mathbf{A}_n = n^{-1/2} \sum_{i=1}^n \dot{\mathbf{l}}(X_i, \boldsymbol{\theta}_0) \xrightarrow{\mathcal{L}} N_k(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}_0))$ .

Next let  $B_{nrs}(\mathbf{t}) = -n^{-1} \sum_{i=1}^n \ddot{l}_{rs}(X_i, \mathbf{t})$ . Then

$$\mathbf{B}_n(\mathbf{t}) = ((B_{nrs}(\mathbf{t}))) = -n^{-1} \sum_{i=1}^n (\ddot{l}_{rs}(X_i, \mathbf{t})) \xrightarrow{P} -\mathbb{E}[(\ddot{l}_{rs}(X, \mathbf{t}))].$$

By Condition 2, arguing in the same manner as in the single-parameter case, we see that  $\mathbf{T}_n = \boldsymbol{\theta}_0 + o_P(1)$  implies  $\mathbf{B}_n(\mathbf{T}_n) = \mathbf{B}_n(\boldsymbol{\theta}_0) + o_P(1)$ . [For vectors and matrices, the  $o_P(1)$  and  $O_P(1)$  notations apply to each coordinate of the vectors and each entry of the matrices.]

By WLLN, we now have  $\mathbf{B}_n(\mathbf{T}_n) = \mathbf{I}(\boldsymbol{\theta}_0) + o_P(1)$  if  $\mathbf{T}_n = \boldsymbol{\theta}_0 + o_P(1)$ .

**Theorem 7.1.3** (Asymptotic Normality of MLEs: Multiparameter Case). *Under regularity Conditions 1, 2, 3, the MLE  $\hat{\boldsymbol{\theta}}_n$  is asymptotically normal (ie,  $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{\mathcal{L}} N_k(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}_0)^{-1})$ ).*

*Proof.* The MLE  $\hat{\boldsymbol{\theta}}_n$  of  $\boldsymbol{\theta}_0$  is the solution of the  $k$  equations

$$\sum_{i=1}^n \dot{l}_r(X_i, \hat{\boldsymbol{\theta}}_n) = 0, \quad r = 1, \dots, k.$$

Expanding, as in the single-parameter case,

$$\begin{aligned} 0 &= n^{-1/2} \sum_{i=1}^n \dot{l}_r(X_i, \hat{\boldsymbol{\theta}}_n) \\ &= n^{-1/2} \sum_{i=1}^n \dot{l}_r(X_i, \boldsymbol{\theta}_0) + n^{-1/2} \sum_{i=1}^n \sum_{s=1}^k (\hat{\theta}_{ns} - \theta_{0s}) \ddot{l}_{rs}(X_i, \boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)) \\ &= n^{-1/2} \sum_{i=1}^n \dot{l}_r(X_i, \boldsymbol{\theta}_0) \\ &\quad - \sum_{s=1}^k n^{1/2} (\hat{\theta}_{ns} - \theta_{0s}) n^{-1} \sum_{i=1}^n \{-\ddot{l}_{rs}(X_i, \boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0))\} \\ &= A_{nr} - \sum_{s=1}^k n^{1/2} (\hat{\theta}_{ns} - \theta_{0s}) B_{nrs}(\boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)), \quad 0 \leq \lambda \leq 1, \\ \text{ie, } &\sum_{s=1}^k n^{1/2} (\hat{\theta}_{ns} - \theta_{0s}) B_{nrs}(\boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)) = A_{nr}, \quad r = 1, \dots, k. \end{aligned}$$

Putting these  $k$  equations together in matrix notation, we get

$$\begin{aligned} \mathbf{B}_n(\boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)) \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) &= \mathbf{A}_n, \text{ and hence} \\ \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) &= \mathbf{B}_n(\boldsymbol{\theta}_0 + \lambda(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0))^{-1} \mathbf{A}_n \\ &= [I(\boldsymbol{\theta}_0) + o_P(1)]^{-1} \mathbf{A}_n \xrightarrow{\mathcal{L}} I(\boldsymbol{\theta}_0)^{-1} \mathbf{W}, \text{ where } \mathbf{W} \sim N_k(\mathbf{0}, I(\boldsymbol{\theta}_0)). \end{aligned}$$

□

Now consider the situation where the likelihood equations do not have a closed form solution. Here again, as in the single-parameter case, we use the Newton-Raphson method, starting with an initial estimator  $\tilde{\theta}_{n0}$  to calculate the next iterate  $\tilde{\theta}_{n1}$  by the formula

$$\tilde{\theta}_{n1} = \tilde{\theta}_{n0} + \mathbf{B}_n(\tilde{\theta}_{n0})^{-1} n^{-1} \sum_{i=1}^n \mathbf{i}(X_i, \tilde{\theta}_{n0}).$$

If the initial estimator  $\tilde{\theta}_{n0}$  is  $\sqrt{n}$ -consistent, then  $\tilde{\theta}_{n1}$  would be a BAN estimator.

*Remark 7.1.1.* We have seen that under some regularity conditions on  $\{f(x, \theta), \theta \in \Theta\}$ , the likelihood equation based on iid observations on  $f(x, \theta_0)$  has a consistent sequence of solutions  $\tilde{\theta}_n$  of  $\theta_0$  if  $\theta_0$  is an interior point of  $\Theta$  and for such a sequence,  $\sqrt{n}(\tilde{\theta}_n - \theta_0) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta_0))$ . However, these are local properties of the likelihood function, and such a sequence of solutions of the likelihood equation need not be the sequence of actual MLEs. On the other hand, globally, under another set of regularity conditions (due to Wald [27]), the actual MLE  $\hat{\theta}_n$  converges almost surely to  $\theta_0$  (strong consistency). There is another aspect of the results about  $\hat{\theta}_n$  discussed above. Due to the fact that the asymptotic variance of  $\tilde{\theta}_n$  is the same as the information lower bound of an unbiased estimator, some sort of asymptotic efficiency is suggested. Is it true that if  $\sqrt{n}(T_n - \theta_0) \xrightarrow{\mathcal{L}} N(0, V(\theta_0))$  for some estimator  $T_n$ , then  $V(\theta_0) \geq 1/I(\theta_0)$ ?

The following examples illustrate how each of these properties of the MLE can be violated. In each of these examples, we describe  $f(x, \theta)$  from which iid observations  $X_1, \dots, X_n$  yield  $\hat{\theta}_n$  or  $\tilde{\theta}_n$ . [Example 7.1.2\(a\)](#) is of a different nature.

**Example 7.1.1.** MLE is not asymptotically normal. Let  $f(x, \theta) = I_{(0,\theta)}(x)$ . Here the MLE  $\hat{\theta}_n = X_{n:n} = \max(X_1, \dots, X_n)$  and

$$P_\theta[n(\theta - X_{n:n}) \leq t] \rightarrow 1 - \exp(-t/\theta) \text{ as } n \rightarrow \infty.$$

(See [Example 7.1.4](#))

**Example 7.1.2** (MLE Is Not Consistent).

**(a) Neyman and Scott [28]**

Let  $(X_i, Y_i)$ ,  $i = 1, 2, \dots$  be independent with  $(X_i, Y_i) \sim N_2((\mu_i, \mu_i), \sigma^2 I)$ . The parameter of interest is  $\sigma^2$  which is to be estimated, while  $\mu_1, \mu_2, \dots$  are nuisance parameters.

The log likelihood based on  $(X_1, Y_1), \dots, (X_n, Y_n)$  is

$$\begin{aligned} \log L &= -n[\log 2\pi + \log \sigma^2] - (2\sigma^2)^{-1} \sum_{i=1}^n [(X_i - \mu_i)^2 + (Y_i - \mu_i)^2] \\ &= -n[\log 2\pi + \log \sigma^2] - (4\sigma^2)^{-1} \sum_{i=1}^n [Z_i^2 + 4\{(X_i + Y_i)/2 - \mu_i\}^2], \end{aligned}$$

where  $Z_i = X_i - Y_i$ . Obviously, the MLE for each  $\mu_i$  is  $\hat{\mu}_i = (X_i + Y_i)/2$  and it follows that the MLE of  $\sigma^2$  is  $\hat{\sigma}_n^2 = (4n)^{-1} \sum_{i=1}^n Z_i^2$ . Now the  $Z_i$ 's are iid  $N(0, 2\sigma^2)$ , so that

$(2n)^{-1} \sum_{i=1}^n Z_i^2$  is a consistent estimator of  $\sigma^2$ . Hence  $\hat{\sigma}_n^2 \xrightarrow{P} \sigma^2/2$ .

**(b) Basu [29]**

Let  $A$  be the set of all rational numbers in  $[0, 1]$  and  $B$  any (known) countable set of irrational numbers in  $[0, 1]$ . The iid observations  $(X_1, \dots, X_n)$  are 0–1 valued rv's which are  $Bernoulli(\theta)$  for  $\theta \in A$  and  $Bernoulli(1 - \theta)$  for  $\theta \in B$ . We therefore find

$$\begin{aligned} M_1 &= \max_{\theta \in A} \theta \sum X_i (1 - \theta)^{n - \sum X_i} = \max_{\theta \in A} L(\theta) \text{ and} \\ M_2 &= \max_{\theta \in B} (1 - \theta) \sum X_i \theta^{n - \sum X_i} = \max_{\theta \in B} L(\theta). \end{aligned}$$

Thus the MLE of  $\theta$  is

$$\hat{\theta}_n = \begin{cases} \arg \max_{\theta \in A} \theta \sum X_i (1 - \theta)^{n - \sum X_i} & \text{if } M_1 > M_2 \\ \arg \max_{\theta \in B} (1 - \theta) \sum X_i \theta^{n - \sum X_i} & \text{if } M_1 < M_2, \end{cases}$$

taking either one if  $M_1 = M_2$ .

Now  $\max_{\theta \in [0, 1]} \theta \sum X_i (1 - \theta)^{n - \sum X_i} = (\sum X_i/n)^{\sum X_i} (1 - \sum X_i/n)^{n - \sum X_i}$  and the maximizer  $\sum X_i/n \in A$  because it is rational. Hence

$$M_1 = \max_{\theta \in A} L(\theta) = \max_{\theta \in [0, 1]} L(\theta) > \max_{\theta \in B} L(\theta) = M_2.$$

Thus  $\hat{\theta}_n = \sum X_i/n$ . But  $\hat{\theta}_n \xrightarrow{P} \theta$  if  $\theta \in A$  and  $\hat{\theta}_n \xrightarrow{P} 1 - \theta$  if  $\theta \in B$ , and is therefore inconsistent. So far we have not used the fact that  $B$  is countable. However, when  $B$  is countable, one can construct a consistent estimator.

**(c) Ferguson [30]**

Let

$$f(x, \theta) = (1 - \theta)f_1(x, \theta) + \theta f_2(x, \theta), \quad \theta \in \Theta = [0, 1],$$

where

$$f_1(x, \theta) = \frac{1}{\delta(\theta)} \left[ 1 - \frac{|x - \theta|}{\delta(\theta)} \right] I_{[\theta - \delta(\theta), \theta + \delta(\theta)]}(x)$$

is a triangular density with base  $[\theta - \delta(\theta), \theta + \delta(\theta)]$  and height  $\delta(\theta)$ , and

$f_2(x, \theta) = (1/2)I_{[-1, 1]}(x)$  is a uniform pdf on  $[-1, 1]$ . The function  $\delta(\theta)$  is a continuous decreasing function of  $\theta$  on  $0 < \theta \leq 1$  with  $\delta(0) = 1$  and  $0 < \delta(\theta) \leq 1 - \theta$ . As  $\theta$  increases from 0 to 1, the triangle's base in  $f_1(x, \theta)$  becomes shorter and shorter, its height becomes larger and larger, and it receives less and less weight. In this way,  $f(x, \theta)$  continuously changes from the triangular to the rectangular density. Now suppose  $\delta(\theta) \rightarrow 0$  as  $\theta \rightarrow 1$  at a sufficiently fast rate so that

$n^{-1} \log[(1 - X_{n:n})/\delta(X_{n:n})] \rightarrow \infty$  with probability 1 for  $\theta = 0$  (triangular case) and hence for all  $\theta$ . Then the MLE  $\hat{\theta}_n \xrightarrow{a.s.} 1$ , whatever the true value  $\theta_0 \in [0, 1]$  may be.

**Remark 7.1.2** (Super Efficiency). The MLE  $\hat{\theta}_n$  has asymptotic variance  $1/I(\theta_0)$ . The possibility of estimators with asymptotic variance  $V(\theta_0) \leq 1/I(\theta_0)$  will be discussed in [Section 7.1.2, Example 7.1.7](#).

**Example 7.1.3.** Let  $(X_1, \dots, X_n)$  be a random sample from  $\text{Cauchy}(\theta, 1)$  with pdf

$$f(x, \theta) = \frac{1}{\pi[1 + (x - \theta)^2]}.$$

Find a BAN estimator of  $\theta$  and its asymptotic distribution.

*Solution.* Here

$$\begin{aligned} l(x, \theta) &= \log f(x, \theta) = -\log \pi - \log[1 + (x - \theta)^2], \\ \dot{l}(x, \theta) &= \frac{2(x - \theta)}{1 + (x - \theta)^2} \text{ and } \ddot{l}(x, \theta) = \frac{2[(x - \theta)^2 - 1]}{[1 + (x - \theta)^2]^2}. \end{aligned}$$

Then

$$\begin{aligned} I(\theta) &= E_\theta[\dot{l}^2(X, \theta)] = \int_{-\infty}^{\infty} \frac{4(x - \theta)^2}{[1 + (x - \theta)^2]^2} \frac{dx}{\pi[1 + (x - \theta)^2]} \\ &= \frac{8}{\pi} \int_0^\infty \frac{y^2 dy}{(1 + y^2)^3} = 1/2. \end{aligned}$$

Start with initial estimator  $\tilde{\theta}_{n0} = \text{Sample Median} = X_{n:[n/2]}$  which is  $\sqrt{n}$ -consistent (as will be shown in [Chapter 8](#)). Then a BAN estimator is

$$\begin{aligned} \tilde{\theta}_{n1} &= \tilde{\theta}_{n0} + \frac{\sum_{i=1}^n \dot{l}(X_i, \tilde{\theta}_{n0})}{\sum_{i=1}^n (-\ddot{l}(X_i, \tilde{\theta}_{n0}))} \\ &= \tilde{\theta}_{n0} + \sum_{i=1}^n \left[ \frac{2(X_i - \tilde{\theta}_{n0})^2}{1 + (X_i - \tilde{\theta}_{n0})^2} \right] \Bigg/ \sum_{i=1}^n \left[ \frac{2[(X_i - \tilde{\theta}_{n0})^2 - 1]}{(1 + (X_i - \tilde{\theta}_{n0})^2)^2} \right], \end{aligned}$$

and  $\sqrt{n}(\tilde{\theta}_{n1} - \theta_0) \xrightarrow{\mathcal{L}} N(0, 2)$ .

**Example 7.1.4.** Let  $(X_1, \dots, X_n)$  be a random sample from  $\text{Unif}(0, \theta)$  with pdf  $f(x, \theta) = \theta^{-1}I_{[0,\theta]}(x), \theta > 0$ . Find the MLE of  $\theta$  and its asymptotic distribution.

*Solution.* Here

$$L(\theta) = \prod_{i=1}^n f(X_i, \theta) = \begin{cases} 0 & \theta < X_{n:n} = \max(X_1, \dots, X_n) \\ \theta^{-n} & \theta > X_{n:n}. \end{cases}$$

It is easy to see that  $L(\theta)$  has a unique maximum at  $X_{n:n}$ . Hence the MLE  $\hat{\theta}_n = X_{n:n}$ . To find the distribution of  $\hat{\theta}_n$ , note that

$$\begin{aligned} P[\hat{\theta}_n \leq t] &= P_\theta[X_i \leq t, i = 1, \dots, n] = \begin{cases} (t/\theta)^n & 0 < t \leq \theta \\ 1 & t > \theta \end{cases}, \text{ and} \\ P[n(\theta - \hat{\theta}_n) \leq t] &= P_\theta[\hat{\theta}_n > \theta - t/n] = 1 - P_\theta[\hat{\theta}_n \leq \theta - t/n] \\ &= 1 - \left( \frac{\theta - t/n}{\theta} \right)^n = 1 - \left( 1 - \frac{t}{n\theta} \right)^n \rightarrow 1 - e^{-t/\theta}. \end{aligned}$$

**Example 7.1.5.** Let  $X_1, \dots, X_n$  be iid  $N_k(\mu, \Sigma)$  where  $\Sigma$  is positive definite. Find the MLE of  $(\mu, \Sigma)$ .

*Solution.* The MLEs for  $\mu$  and  $\Sigma$  are  $\bar{X} = n^{-1} \sum_{i=1}^n X_i$  and  $\hat{\Sigma} = n^{-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$ . A proof is given later in the chapter on multivariate analysis.

**Example 7.1.6.** Let

$$f(x, \theta) = C(\theta) \exp[Q(\theta)T(x)]r(x), \quad \theta \in \Theta$$

be a pdf/pmf belonging to the single-parameter exponential family, where  $\Theta = \{\theta : \int \exp[Q(\theta)T(x)]r(x) dx < \infty\}$ . Find the MLE of  $\theta$  from a random sample  $(X_1, \dots, X_n)$  from  $f(x, \theta)$ .

*Solution.* The log-likelihood and the likelihood equation are

$$\begin{aligned} l(\mathbf{X}, \theta) &= n \log C(\theta) + Q(\theta) \sum_{i=1}^n T(X_i) + \sum_{i=1}^n \log r(X_i), \text{ and} \\ \dot{l}(\mathbf{X}, \theta) &= n(C'(\theta)/C(\theta)) + Q'(\theta) \sum_{i=1}^n T(X_i) = 0. \end{aligned}$$

The MLE  $\hat{\theta}_n$  is the solution of

$$-\frac{C'(\theta)}{C(\theta)Q'(\theta)} = n^{-1} \sum_{i=1}^n T(X_i).$$

We can now find the MLEs of  $\theta$  for some distributions as special cases of this result.

**(a) Bernoulli( $\theta$ )** with

$$f(x, \theta) = \theta^x (1 - \theta)^{1-x} = (1 - \theta) \exp\left[x \log\left(\frac{\theta}{1 - \theta}\right)\right], \quad x = 0, 1.$$

Here

$$C(\theta) = 1 - \theta, \quad Q(\theta) = \log\left(\frac{\theta}{1 - \theta}\right), \quad \text{and} \quad T(x) = x, \quad \text{so} \quad -\frac{C'(\theta)}{C(\theta)Q'(\theta)} = \theta,$$

and the likelihood equation is  $\theta = n^{-1} \sum_{i=1}^n X_i = \bar{X}$ . Thus the MLE of  $\theta$  is  $\hat{\theta}_n = \bar{X}$ .

**(b) Poisson( $\theta$ )** with

$$f(x, \theta) = \exp(-\theta) \theta^x / x!, \quad x = 0, 1, \dots$$

Here

$$C(\theta) = e^{-\theta}, \quad Q(\theta) = \log \theta, \quad T(x) = x, \quad \text{so} \quad -\frac{C'(\theta)}{C(\theta)Q'(\theta)} = \theta,$$

and the likelihood equation is  $\theta = n^{-1} \sum_{i=1}^n X_i = \bar{X}$ . Thus the MLE of  $\theta$  is  $\hat{\theta}_n = \bar{X}$ .

**(c) Gamma( $\alpha, \theta$ )** with

$$f(x, \theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha} \exp[-x/\theta] x^{\alpha-1}, \quad x > 0.$$

Here

$$C(\theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha}, \quad Q(\theta) = -1/\theta, \quad T(x) = x, \quad \text{so} \quad -\frac{C'(\theta)}{C(\theta)Q'(\theta)} = \alpha\theta,$$

and the likelihood equation is  $\alpha\theta = n^{-1} \sum_{i=1}^n X_i = \bar{X}$ . Thus the MLE of  $\theta$  is  $\hat{\theta}_n = \bar{X}/\alpha$ .

### 7.1.1 Almost Sure Convergence (Strong Consistency) of MLEs

Let  $\theta_0$  be the true value of the unknown parameter  $\theta$  and all probability statements are with respect to  $\theta_0$ .

#### Notations

- (i)** For every  $\theta \in \Theta$ , let  $N_{\theta,j} = \{\varphi: \|\varphi - \theta\| \leq j^{-1}\}$  be a decreasing sequence of neighborhoods of  $\theta$  converging to  $\{\theta\}$ , and let

$$Z(\theta_0, N_{\theta,j}) = \inf_{\varphi \in N_{\theta,j}} \log[f(X, \theta_0)/f(X, \varphi)], \quad Z(\theta_0, \theta) = \log[f(X, \theta_0)/f(X, \theta)].$$

- (ii)** Let  $N_{\infty,j} = \{\theta: \|\theta\| > j\}$ , which can be thought of as a decreasing sequence of neighborhoods of  $\infty$ , converging to  $\emptyset$ , and let

$$Z(\theta_0, N_{\infty,j}) = \inf_{\varphi \in N_{\infty,j}} [f(X, \theta_0)/f(X, \varphi)].$$

We assume that the following *Regularity Conditions (Wald)* hold

1. The parameter space  $\Theta$  is a closed subset of  $\mathbb{R}^k$ .
2. *Identifiability Condition.* For  $\theta \neq \theta_0$ ,  $\{x: f(x, \theta) \neq f(x, \theta_0)\}$  has positive probability under  $\theta_0$ .
3. **(a)** For all  $\theta$ ,  $\lim_{\theta_n \rightarrow \theta} f(x, \theta_n) = f(x, \theta)$  for all  $x$ .  
**(b)**  $\lim_{\theta_n \rightarrow \infty} f(x, \theta_n) = 0$  for all  $x$ .
4. **(a)** For each  $\theta$ ,  $I(\theta_0, \theta) = E[Z(\theta_0, \theta)]$  exists.  
**(b)** For each  $\theta$ ,  $I(\theta_0, N_{\theta,j}) = E[Z(\theta_0, N_{\theta,j})] > -\infty$  for some  $j = j_0$ .  
**(c)**  $I(\theta_0, N_{\infty,j}) = E[Z(\theta_0, N_{\infty,j})] > -\infty$  for some  $j = j_0$ .

**Theorem 7.1.4** (Wald). *Under Conditions 1–4,  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$ .*

The proof will be accomplished by first considering the case of finite  $\Theta = \{\theta_0, \theta_1, \dots, \theta_r\}$  and then extending to the case of  $\Theta$  being an arbitrary closed subset of  $\mathbb{R}^k$ .

We start with the following lemma.

**Lemma 7.1.1.**  $I(\theta_0, \theta) > 0$  for all  $\theta \neq \theta_0$ .

This result has already been proved in [Section 5.5.1](#), using the Identifiability Condition.

In the case of finite  $\Theta$ , we now have

$$\begin{aligned} P[\hat{\theta}_n \neq \theta_0 \text{ i.o.}] &\leq \sum_{j=1}^r P[\hat{\theta}_n = \theta_j \text{ i.o.}] \\ &\leq \sum_{j=1}^r P\left[n^{-1} \sum_{i=1}^n \log f(X_i, \theta_0) \leq n^{-1} \sum_{i=1}^n \log f(X_i, \theta_j) \text{ i.o.}\right] \\ &= \sum_{j=1}^r P\left[n^{-1} \sum_{i=1}^n Z_i(\theta_0, \theta_j) \leq 0 \text{ i.o.}\right] = 0 \end{aligned}$$

by SLLN, because  $E[Z(\theta_0, \theta_j)] = I(\theta_0, \theta_j) > 0$  for all  $j$  by Lemma 7.1.1. This proves the theorem for finite  $\Theta$ .

For the case of  $\Theta$  being an arbitrary closed subset of  $\mathbb{R}^k$ , we need to show that for an arbitrary neighborhood  $N_0$  of  $\theta_0$ ,  $P[\hat{\theta}_n \notin N_0 \text{ i.o.}] = 0$ . To extend the proof for the finite case to this generality, we shall cover  $\Theta - N_0$  by a finite collection of sets  $S_1, \dots, S_r$  such that

$$I(\theta_0, S_j) = E[Z(\theta_0, S_j)] = E\left[\inf_{\varphi \in S_j} \log\{f(X, \theta_0)/f(X, \varphi)\}\right] > 0$$

for  $j = 1, \dots, r$ . It then follows that

$$\begin{aligned} P[\hat{\theta}_n \notin N_0 \text{ i.o.}] &\leq P[\hat{\theta}_n \in S_j \text{ i.o. for some } j = 1, \dots, r] \\ &\leq \sum_{j=1}^r P[\hat{\theta}_n \in S_j \text{ i.o.}] \\ &\leq \sum_{j=1}^r P\left[n^{-1} \sum_{i=1}^n \log f(X_i, \theta_0) \leq \sup_{\varphi \in S_j} n^{-1} \sum_{i=1}^n \log f(X_i, \varphi) \text{ i.o.}\right] \\ &= \sum_{j=1}^r P\left[\inf_{\varphi \in S_j} n^{-1} \sum_{i=1}^n \log\{f(X_i, \theta_0)/f(X_i, \varphi)\} \leq 0 \text{ i.o.}\right] \\ &\leq \sum_{j=1}^r P\left[n^{-1} \sum_{i=1}^n \inf_{\varphi \in S_j} \log\{f(X_i, \theta_0)/f(X_i, \varphi)\} \leq 0 \text{ i.o.}\right] \\ &= \sum_{j=1}^r P\left[n^{-1} \sum_{i=1}^n Z_i(\theta_0, S_j) \leq 0 \text{ i.o.}\right] = 0 \end{aligned}$$

by SLLN because  $E[Z(\theta_0, S_j)] > 0$  for  $j = 1, \dots, r$ .

We now construct  $S_1, \dots, S_r$  with the above-mentioned properties.

**Lemma 7.1.2.**

- (a) For each  $\theta \neq \theta_0$ , there is a neighborhood  $N_{\theta,j}$  such that  $I(\theta_0, N_{\theta,j}) = E[Z(\theta_0, N_{\theta,j})] > 0$ .
- (b) There exists a neighborhood  $N_{\infty,j}$  of  $\infty$  such that  $I(\theta_0, N_{\infty,j}) = E[Z(\theta_0, N_{\infty,j})] > 0$ .

*Proof.* Fix  $\theta \neq \theta_0$  and let

$$\begin{aligned} g_j(x) &= \left\{ \inf_{\varphi \in N_{\theta,j}} \log(f(x, \theta_0)/f(x, \varphi)) \right\} f(x, \theta_0) \\ &\quad - \left\{ \inf_{\varphi \in N_{\theta,j_0}} \log(f(x, \theta_0)/f(x, \varphi)) \right\} f(x, \theta_0) \end{aligned}$$

for  $j \geq j_0$ , where  $j_0$  is chosen by Condition 4(b). Since  $\{N_{\theta,j}\}$  is decreasing,  $\{g_j\}$  is increasing and are clearly nonnegative for  $j \geq j_0$ . Moreover, by Condition 3(a),

$$\begin{aligned} \lim_{j \rightarrow \infty} g_j(x) &:= g(x) = \{\log(f(x, \theta_0)/f(x, \theta))\} f(x, \theta_0) \\ &\quad - \left\{ \inf_{\varphi \in N_{\theta,j_0}} \log(f(x, \theta_0)/f(x, \varphi)) \right\} f(x, \theta_0) \end{aligned}$$

for all  $x$ . Hence by the Monotone Convergence Theorem,

$$\begin{aligned} \lim_{j \rightarrow \infty} E[Z(\theta_0, N_{\theta,j})] - E[Z(\theta_0, N_{\theta,j_0})] &= \lim_{j \rightarrow \infty} \int g_j(x) dx = \int g(x) dx \\ &= E[Z(\theta_0, \theta)] - E[Z(\theta_0, N_{\theta,j_0})]. \end{aligned}$$

Since  $E[Z(\theta_0, N_{\theta,j_0})] > -\infty$ , we cancel it from both sides of the above equality to obtain

$$\lim_{j \rightarrow \infty} E[Z(\theta_0, N_{\theta,j})] = E[Z(\theta_0, \theta)] > 0,$$

so that  $E[Z(\theta_0, N_{\theta,j})] > 0$  for sufficiently large  $j$ . This proves (a).

The proof of (b) is exactly the same, using Conditions 4(c) and 3(b) instead of Conditions 4(b) and 3(a).  $\square$

*Proof of Theorem 7.1.4.* First choose a set  $N_{\infty,j^*}$  such that  $E[Z(\theta_0, N_{\infty,j^*})] > 0$ . Now consider  $\Theta \cap N_0^c \cap N_{\infty,j^*}^c$  (for an arbitrary neighborhood  $N_0$  of  $\theta_0$ ), which is a closed and bounded set in  $\mathbb{R}^k$  and therefore has the Heine-Borel property (compactness) of having a subcover for every open cover.

Now for every  $\theta \in \Theta \cap N_0^c \cap N_{\infty,j^*}^c$ , choose  $N_{\theta,j}$  such that  $E[Z(\theta_0, N_{\theta,j})] > 0$ . These sets  $N_{\theta,j}$  for  $\theta \in \Theta \cap N_0^c \cap N_{\infty,j^*}^c$  form an open cover of  $\Theta \cap N_0^c \cap N_{\infty,j^*}^c$  and by the Heine-Borel property, there is a finite subcover  $N_{\theta_1,j_1}, \dots, N_{\theta_{r-1},j_{r-1}}$ . Let

$$S_1 = N_{\theta_1,j_1}, \dots, S_{r-1} = N_{\theta_{r-1},j_{r-1}} \text{ and } S_r = N_{\infty,j^*}.$$

Then  $S_1, \dots, S_r$  cover  $N_0^c$  and  $E[Z(\theta_0, S_j)] > 0$ ,  $j = 1, \dots, r$ , and the theorem is proved by the argument given above.  $\square$

### 7.1.2 Asymptotic Efficiency of MLE

What is the connection between the information bound  $1/I(\theta) = 1/E_\theta[\hat{I}^2(X)]$  for the variance of an unbiased estimator of  $\theta$  and the fact that the MLE of  $\hat{\theta}_n$  of  $\theta$  has the

asymptotic distribution  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta))$ ? [For simplicity, we are considering the case of a 1-dim parameter  $\theta$ .]

We have seen that

- (i) under regularity conditions on the pdf/pmf and restrictions on estimators,  $n\text{Var}_\theta[T_n] \geq 1/I(\theta)$  for all unbiased estimators  $T_n$  of  $\theta$ , and
- (ii) the MLE  $\hat{\theta}_n$  of  $\theta$  has the property:  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} N(0, 1/I(\theta))$ .

From (i) and (ii), one may be tempted to hope that if  $\sqrt{n}(T_n - \theta) \xrightarrow{\mathcal{L}} \xi_\theta$  under  $f(\cdot, \theta)$ , then  $\xi_\theta$  must be at least as much dispersed as  $N(0, 1/I(\theta))$  and in particular,  $\sigma^2(\theta) = E_\theta[\xi_\theta^2] \geq 1/I(\theta)$  if  $\xi_\theta$  is Gaussian. Can there be a situation in which  $E_\theta[\xi_\theta^2] < 1/I(\theta)$ ? If so, then such an estimator would be called “superefficient.” The following example shows that *superefficient estimators do exist*.

**Example 7.1.7** (Hodges). Let  $X_1, \dots, X_n, \dots$  be iid as  $N(\theta, 1)$ . Then the MLE  $\hat{\theta}_n = \bar{X}_n = n^{-1} \sum_{i=1}^n X_i$  and  $Z = \sqrt{n}(\bar{X}_n - \theta) \xrightarrow{\mathcal{D}} N(0, 1)$ . Also the information  $I(\theta) = 1$ . Now for  $0 < a < 1$ , let

$$T_n = \begin{cases} \bar{X}_n & \text{if } |\bar{X}_n| > n^{-1/4} \\ a\bar{X}_n & \text{if } |\bar{X}_n| \leq n^{-1/4}. \end{cases}$$

Then  $\sqrt{n}(T_n - \theta)$  is distributed as  $Z$  if  $|Z + \sqrt{n}\theta| > n^{1/4}$  and as  $aZ + (a-1)\sqrt{n}\theta$  otherwise. From this it follows that  $\sqrt{n}(T_n - \theta) \xrightarrow{\mathcal{L}} N(0, V(\theta))$ , where  $V(\theta) = 1$  if  $\theta \neq 0$  and  $V(\theta) = a^2 < 1$  if  $\theta = 0$ . This defines a superefficient estimator (see [31]).

Let us examine the behavior of  $\{T_n\}$  in the above example under  $\{P_{\theta_n}\}$  where  $\theta_n = h/\sqrt{n}$ . In the above proof, taking  $\theta = \theta_n$ , we have

$$\begin{aligned} \sqrt{n}(T_n - \theta_n) &\xrightarrow{\mathcal{D}} ZI[|Z + h| > n^{1/4}] + \{aZ + (a-1)h\}I[|Z + h| \leq n^{1/4}] \\ &\xrightarrow{a.s.} aZ + (a-1)h. \end{aligned}$$

Hence under  $\{P_{\theta_n}\}$ ,  $\sqrt{n}(T_n - \theta_n) \xrightarrow{\mathcal{L}} N((a-1)h, a^2)$ , although for  $\theta = 0$ ,  $\sqrt{n}(T_n - 0) \xrightarrow{\mathcal{L}} N(0, a^2)$ . Thus we have a situation in which the asymptotic distribution of  $\sqrt{n}(T_n - \theta_n)$  under  $\{P_{\theta_n}\}$  with  $\theta_n = h/\sqrt{n} \rightarrow 0$  as  $n \rightarrow \infty$ , depends not only on  $P_0$  but actually on the particular sequence  $\{P_{\theta_n}\} \rightarrow P_0$ .

Hájek recognized this as the reason behind superefficiency and introduced the concept of *regular estimators* (or more precisely, locally regular estimators at a particular  $\theta_0$ ) to prevent superefficiency.

**Definition 7.1.2.** A sequence of estimators  $\{T_n\}$  of  $\theta$  is regular at  $P = P_{\theta_0}$  if for every sequence  $\{P_n\} = \{P_{\theta_n}\}$  with  $\theta_n = \theta_0 + h/\sqrt{n} + o(1/\sqrt{n})$  for some  $h \in \mathbb{R}$ ,

$$\mathcal{L}(\sqrt{n}(T_n - \theta_n)) \rightarrow \mathcal{L}(P) \text{ under } \{P_n\},$$

where  $\mathcal{L}(P)$  depends on  $P = P_{\theta_0}$ , but not on the particular sequence  $\{P_{\theta_n}\}$  (ie, not on  $h$ ).

**Theorem 7.1.5** (Hájek-Inagaki Decomposition Theorem) [See 67, 68]. Suppose that  $\{T_n\}$  is a regular estimator of  $\theta$  at  $P = P_{\theta_0}$ . Then for every sequence  $\theta_n = \theta_0 + h/\sqrt{n} + o(1/\sqrt{n})$  with  $h \in \mathbb{R}$ ,

$$\mathcal{L}(\sqrt{n}(T_n - \theta_n)) \rightarrow \mathcal{L}(P) = N(0, 1/I) * \mathcal{L}_1(P) \text{ under } P_n,$$

where  $\mathcal{L}_1(P)$  depends only on  $P$  (ie, on  $\theta_0$ ), but not on the particular sequence  $\{\theta_n\}$ , and  $*$  denotes convolution.

By this theorem, the asymptotic distribution of any regular estimator is more dispersed than  $N(0, 1/I)$ , which is the asymptotic distribution of the MLE. Thus *the MLE is asymptotically efficient among all regular estimators*.

### 7.1.3 MLE of Parameters in a Multinomial Distribution

Consider a multinomial distribution in  $m$  classes with probability  $\pi_j(\theta)$ ,  $j = 1, \dots, m$  for the  $j$ th class, where  $\pi_1(\cdot), \dots, \pi_m(\cdot)$  are known functions of an unknown  $k$ -dim parameter vector  $\theta$ . Let  $n_1, \dots, n_m$  denote the observed frequencies in the  $m$  classes in a random sample of size  $n$ . The following table summarizes this.

Class	1	...	$j$	...	$m$	Total
Probability	$\pi_1(\theta)$	...	$\pi_j(\theta)$	...	$\pi_m(\theta)$	1
Obs. frequency	$n_1$	...	$n_j$	...	$n_m$	$n$

Here the data consist of a random sample  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  where the  $\mathbf{X}_i$ s are iid with

$$P[\mathbf{X}_i = \mathbf{e}_j] = f(\mathbf{e}_j, \theta) = \pi_j(\theta), \quad j = 1, \dots, m \text{ and } n_j = \sum_{i=1}^n I[\mathbf{X}_i = \mathbf{e}_j],$$

$\mathbf{e}_j$  being the  $m$ -dim  $j$ th coordinate vector. Hence

$$\begin{aligned} l_r(\mathbf{e}_j, \theta) &= \frac{\partial}{\partial \theta_r} \log f(\mathbf{e}_j, \theta) = \frac{\partial}{\partial \theta_r} \log \pi_j(\theta) = \frac{1}{\pi_j(\theta)} \frac{\partial \pi_j(\theta)}{\partial \theta_r} := \frac{\dot{\pi}_{rj}(\theta)}{\pi_j(\theta)}, \text{ and} \\ S_{nr}(\mathbf{t}) &= \sum_{i=1}^n l_r(\mathbf{X}_i, \mathbf{t}) = \sum_{i=1}^n \sum_{j=1}^m I[\mathbf{X}_i = \mathbf{e}_j] \frac{\dot{\pi}_{rj}(\mathbf{t})}{\pi_j(\mathbf{t})} = \sum_{j=1}^m \frac{n_j}{\pi_j(\mathbf{t})} \dot{\pi}_{rj}(\mathbf{t}). \end{aligned}$$

Thus the likelihood equations are

$$S_{nr}(\hat{\theta}_n) = 0, \text{ ie, } \sum_{j=1}^m \frac{n_j}{\pi_j(\hat{\theta}_n)} \dot{\pi}_{rj}(\hat{\theta}_n) = 0, \quad r = 1, \dots, k.$$

In a typical context, these equations would not have a closed form solution, so we obtain a BAN estimator by the Newton-Raphson method. First note that

$$P_\theta \left[ l_r(\mathbf{X}, \theta) = \frac{\dot{\pi}_{rj}(\theta)}{\pi_j(\theta)} \right] = P_\theta[\mathbf{X} = \mathbf{e}_j] = \pi_j(\theta), \quad j = 1, \dots, m.$$

For  $\theta = \theta_0$ ,

$$\mathbb{E}_{\theta_0} [\dot{l}_r(\mathbf{X}, \theta_0)] = \sum_{j=1}^m \pi_j(\theta_0) \frac{\dot{\pi}_{rj}(\theta_0)}{\pi_j(\theta_0)} = \sum_{j=1}^m \dot{\pi}_{rj}(\theta_0) = \left. \frac{\partial}{\partial \theta_r} \sum_{j=1}^m \pi_j(\theta) \right|_{\theta=\theta_0} = 0,$$

as it should be because  $\sum_{j=1}^m \pi_j(\theta) = 1$  for all  $\theta$ . Next,

$$\begin{aligned} I_{rs}(\theta_0) &= \text{Cov}_{\theta_0} [\dot{l}_r(\mathbf{X}, \theta_0), \dot{l}_s(\mathbf{X}, \theta_0)] = \mathbb{E}_{\theta_0} [\dot{l}_r(\mathbf{X}, \theta_0) \dot{l}_s(\mathbf{X}, \theta_0)] \\ &= \sum_{j=1}^m \pi_j(\theta_0) \{\dot{\pi}_{rj}(\theta_0)/\pi_j(\theta_0)\} \{\dot{\pi}_{sj}(\theta_0)/\pi_j(\theta_0)\} \\ &= \sum_{j=1}^m \dot{\pi}_{rj}(\theta_0) \dot{\pi}_{sj}(\theta_0)/\pi_j(\theta_0). \end{aligned} \quad (5)$$

If  $\mathbf{T}_n = \theta_0 + o_P(1)$ , then  $I_{rs}(\theta_0)$  is estimated by  $I_{rs}(\mathbf{T}_n)$ , which can be used in place of  $B_{nrs}(\tilde{\theta}_n)$  in the Newton-Raphson formula.

## 7.2 Likelihood Ratio Test

The theory of hypothesis testing, discussed in [Chapter 6](#), based on N-P Lemma for simple  $H_0$  vs simple  $H_1$  and extended to some composite  $H_0$  vs composite  $H_1$  in exponential families, also in the presence of nuisance parameters, using the Generalized N-P Lemma, covered a limited range of problems. For testing a composite  $H_0: \theta \in \Theta_0$  vs a composite  $H_1: \theta \in \Theta_1$  in a general setting, we modify the likelihood ratio  $\prod_{i=1}^n f(X_i, \theta_1)/\prod_{i=1}^n f(X_i, \theta_0)$  based on  $\mathbf{X} = (X_1, \dots, X_n)$  in the simple vs simple problem by replacing  $\theta_0$  and  $\theta_1$  by  $\hat{\theta}_i = \text{MLE of } \theta \text{ restricted to } \Theta_i$ ,  $i = 0, 1$ , respectively, or equivalently by

$$\Lambda_n = \frac{\sup_{\theta \in \Theta_0} \prod_{i=1}^n f(X_i, \theta)}{\sup_{\theta \in \Theta} \prod_{i=1}^n f(X_i, \theta)} = \frac{f(X_i, \tilde{\theta}_n)}{f(X_i, \hat{\theta}_n)}, \quad (6a)$$

where  $\tilde{\theta}_n$  is the restricted MLE in  $\Theta_0$  and  $\hat{\theta}_n$  is the unrestricted MLE in the entire parameter space  $\Theta = \mathbb{R}^k$ . The null hypothesis  $H_0$  is rejected for small values of  $\Lambda_n$ , or equivalently, for large values of

$$-2 \log \Lambda_n = 2 \left[ \sum \log f(X_i, \hat{\theta}_n) - \sum \log f(X_i, \tilde{\theta}_n) \right]. \quad (6b)$$

This is known as the likelihood ratio test (LRT). The asymptotic properties of LRT will be discussed in this section.

Hypothesis testing was discussed in [Chapter 6](#), mostly within the very restricted framework of exponential families of distributions. Here we take up the testing problem again in much broader context. Let  $(X_1, \dots, X_n)$  be a random sample from  $f(x, \theta)$ ,  $\theta \in \mathbb{R}^k$ .

We now consider the problem of testing  $H_0: \theta \in \Theta_0$  vs  $H_1: \theta \in \mathbb{R}^k - \Theta_0$ , where  $\Theta_0$  is a  $d$ -dim hyperplane in  $\mathbb{R}^k$  with  $d < k$ . Without loss of generality, let  $\Theta_0 = \{\theta = (\theta_1, \dots, \theta_k): \theta_{d+1} = \dots = \theta_k = 0\}$ .

The Neyman-Pearson simple-vs-simple likelihood ratio would generalize to the composite-vs-composite case as

$$\Lambda_n^* = \frac{\sup_{\theta \in \Theta_0} \prod_{i=1}^n f(X_i, \theta)}{\sup_{\theta \in \mathbb{R}^k - \Theta_0} \prod_{i=1}^n f(X_i, \theta)},$$

and reject  $H_0$  if  $\Lambda_n^* \leq$  critical value  $c$ . However, since determining the MLE of  $\theta$  restricted to  $\mathbb{R}^k - \Theta_0$  causes complication, it is more convenient to modify the denominator and use the test statistic

$$\Lambda_n = \frac{\sup_{\theta \in \Theta_0} \prod_{i=1}^n f(X_i, \theta)}{\sup_{\theta \in \mathbb{R}^k} \prod_{i=1}^n f(X_i, \theta)} = \frac{\prod_{i=1}^n f(X_i, \tilde{\theta}_n)}{\prod_{i=1}^n f(X_i, \hat{\theta}_n)} \text{ as in Eq. (6a),}$$

where  $\tilde{\theta}_n$  is the MLE of  $\theta$  restricted to  $\Theta_0$ ,  $\hat{\theta}_n$  is the unrestricted MLE of  $\theta$  based on  $(X_1, \dots, X_n)$  and  $H_0$  is rejected if  $\Lambda_n \leq$  critical value  $c$ . Clearly,  $\Lambda_n = \Lambda_n^* < 1$  if  $\theta_n \in \mathbb{R}^k - \Theta_0$  and  $\Lambda_n = 1 \leq \Lambda_n^*$  if  $\hat{\theta}_n \in \Theta_0$ . Thus  $\Lambda_n$  is a nondecreasing function of  $\Lambda_n^*$  and therefore, rejecting  $H_0$  for small values of  $\Lambda_n$  is equivalent to rejecting  $H_0$  for small values of  $\Lambda_n^*$ .

We are thus led to the LRT: Reject  $H_0$  for *small values of  $\Lambda_n$*  defined above, or equivalently, reject  $H_0$  for *large values of*

$$-2 \log \Lambda_n = 2 \left[ \sum_{i=1}^n l(X_i, \hat{\theta}_n) - \sum_{i=1}^n l(X_i, \tilde{\theta}_n) \right] \text{ as in Eq. (6b).}$$

To construct a large-sample test of  $H_0$  at a prescribed level of significance, we need the asymptotic distribution of  $-2 \log \Lambda_n$  under  $H_0$ . Without loss of generality, let the true parameter  $\theta_0 = 0$  (ie,  $\theta_{01} = \dots = \theta_{0k} = 0$ ). Since  $\tilde{\theta}_n$  is the MLE of  $\theta$  restricted to  $\Theta_0$ , we must have  $\tilde{\theta}_n = (\tilde{\theta}_{n1}, \dots, \tilde{\theta}_{nd}, 0, \dots, 0)$ .

Thus  $\hat{\theta}_n - \theta_0 = (\hat{\theta}_{n1}, \dots, \hat{\theta}_{nd}, \hat{\theta}_{n,d+1}, \dots, \hat{\theta}_{nk})$  and  $\tilde{\theta}_n - \theta_0 = (\tilde{\theta}_{n1}, \dots, \tilde{\theta}_{nd}, 0, \dots, 0)$ .

We use the notations of Section 5.2.3 and assume that the regularity conditions introduced for proving Theorem 7.1.2 hold. We also write

$$\begin{aligned} \dot{l}_{(d)}(x, \theta)^T &= \dot{l}_1(x, \theta), \dots, \dot{l}_d(x, \theta), \tilde{\theta}_{n(d)}^T = (\tilde{\theta}_{n1}, \dots, \tilde{\theta}_{nd}), \\ \ddot{l}_{(d)}(x, \theta) &= \begin{bmatrix} \ddot{l}_{11}(x, \theta) & \dots & \ddot{l}_{1d}(x, \theta) \\ \vdots & & \vdots \\ \ddot{l}_{d1}(x, \theta) & \dots & \ddot{l}_{dd}(x, \theta) \end{bmatrix} \text{ and } \mathbf{I}_{(d)} = \begin{bmatrix} I_{11} & \dots & I_{1d} \\ \vdots & & \vdots \\ I_{d1} & \dots & I_{dd} \end{bmatrix}, \end{aligned}$$

that is,  $\dot{l}_{(d)}(x, \theta)$  and  $\tilde{\theta}_{n(d)}$  are  $d$ -dim vectors consisting of the first  $d$  elements of  $\dot{l}(x, \theta)$  and  $\tilde{\theta}_n$ , respectively, and  $\ddot{l}_{(d)}(x, \theta)$  and  $\mathbf{I}_{(d)}$  are  $d \times d$  matrices consisting of the upper left-hand elements of  $\ddot{l}(x, \theta)$  and  $\mathbf{I}$ , respectively.

**Theorem 7.2.1.** Under regularity Conditions 1, 2, 3 (multiparameter),  $-2 \log \Lambda_n \xrightarrow{\mathcal{L}} \chi_{k-d}^2$ .

*Proof.* Expanding  $\sum_{i=1}^n l(X_i, \hat{\theta}_n)$  and  $\sum_{i=1}^n l(X_i, \tilde{\theta}_n)$  around  $\theta_0 = \mathbf{0}$ , we have

$$\begin{aligned}\sum_{i=1}^n l(X_i, \hat{\theta}_n) &= \sum_{i=1}^n l(X_i, \mathbf{0}) + (\sqrt{n}\hat{\theta}_n)^T n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \mathbf{0}) \\ &\quad + (1/2)(\sqrt{n}\hat{\theta}_n)^T \left\{ n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \mathbf{0}) + \mathbf{R}_n \right\} (\sqrt{n}\hat{\theta}_n) \\ &= \sum_{i=1}^n l(X_i, \mathbf{0}) + (\sqrt{n}\hat{\theta}_n)^T n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \mathbf{0}) - (1/2)(\sqrt{n}\hat{\theta}_n)^T \mathbf{I}(\mathbf{0})(\sqrt{n}\hat{\theta}_n) + o_P(1),\end{aligned}$$

treating the remainder term  $\mathbf{R}_n$ , using Condition 2 as in [Section 7.1.1](#).

Similarly,

$$\begin{aligned}\sum_{i=1}^n l(X_i, \tilde{\theta}_n) &= \sum_{i=1}^n l(X_i, \mathbf{0}) + (\sqrt{n}\tilde{\theta}_{n(d)})^T n^{-1/2} \sum_{i=1}^n \dot{l}_{(d)}(X_i, \mathbf{0}) \\ &\quad + (1/2)(\sqrt{n}\tilde{\theta}_{n(d)})^T \left\{ n^{-1} \sum_{i=1}^n \ddot{l}_{(d)}(X_i, \mathbf{0}) + \mathbf{R}_n \right\} (\sqrt{n}\tilde{\theta}_{n(d)}) \\ &= \sum_{i=1}^n l(X_i, \mathbf{0}) + (\sqrt{n}\tilde{\theta}_{n(d)})^T n^{-1/2} \sum_{i=1}^n \dot{l}_{(d)}(X_i, \mathbf{0}) \\ &\quad - (1/2)(\sqrt{n}\tilde{\theta}_{n(d)})^T \mathbf{I}_{(d)}(\mathbf{0})(\sqrt{n}\tilde{\theta}_{n(d)}) + o_P(1).\end{aligned}$$

In the above expression, we now substitute (as seen in the proof of [Theorem 7.1.3](#)),

$$\begin{aligned}\sqrt{n}\hat{\theta}_n &= \sqrt{n}(\hat{\theta}_n - \theta_0) = [\mathbf{I}(\mathbf{0}) + o_P(1)]^{-1} n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \mathbf{0}) + o_P(1), \text{ and similarly} \\ \sqrt{n}\tilde{\theta}_{n(d)} &= \sqrt{n}(\tilde{\theta}_{n(d)} - \theta_0) = [\mathbf{I}_{(d)}(\mathbf{0}) + o_P(1)]^{-1} n^{-1/2} \sum_{i=1}^n \dot{l}_{(d)}(X_i, \mathbf{0}) + o_P(1)\end{aligned}$$

to obtain after simplification,

$$\begin{aligned}-2 \log \Lambda_n &= 2 \left[ \sum_{i=1}^n l(X_i, \hat{\theta}_n) - \sum_{i=1}^n l(X_i, \tilde{\theta}_n) \right] \\ &= \left\{ n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \mathbf{0}) \right\}^T \mathbf{I}(\mathbf{0})^{-1} \left\{ n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \mathbf{0}) \right\} \\ &\quad - \left\{ n^{-1/2} \sum_{i=1}^n \dot{l}_{(d)}(X_i, \mathbf{0}) \right\}^T \mathbf{I}_{(d)}(\mathbf{0})^{-1} \left\{ n^{-1/2} \sum_{i=1}^n \dot{l}_{(d)}(X_i, \mathbf{0}) \right\} + o_P(1).\end{aligned}$$

By multivariate CLT,

$$n^{-1/2} \sum_{i=1}^n \mathbf{i}(X_i, \mathbf{0}) \xrightarrow{\mathcal{L}} \mathbf{Y} \sim N_k(\mathbf{0}, \mathbf{I}(\mathbf{0})) \text{ and } n^{-1/2} \sum_{i=1}^n \mathbf{i}_{(d)}(X_i, \mathbf{0}) \xrightarrow{\mathcal{L}} \mathbf{Y}_{(d)},$$

where  $\mathbf{Y}_{(d)}$  is the  $d$ -dim random vector consisting of the first  $d$  coordinates of  $\mathbf{Y}$ . Hence

$$-2 \log \Lambda_n \xrightarrow{\mathcal{L}} \mathbf{Y}^T \mathbf{I}(\mathbf{0})^{-1} \mathbf{Y} - \mathbf{Y}_{(d)}^T \mathbf{I}_{(d)}(\mathbf{0})^{-1} \mathbf{Y}_{(d)}.$$

Using a suitable linear transformation,  $\mathbf{Y}^T \mathbf{I}(\mathbf{0})^{-1} \mathbf{Y}$  and  $\mathbf{Y}_{(d)}^T \mathbf{I}_{(d)}(\mathbf{0})^{-1} \mathbf{Y}_{(d)}$  can be simultaneously reduced to  $\sum_{i=1}^k Z_i^2$  and  $\sum_{i=1}^d Z_i^2$ , respectively, where  $Z_1, \dots, Z_k$  are iid  $N(0, 1)$ . Thus

$$-2 \log \Lambda_n \xrightarrow{\mathcal{L}} \sum_{i=1}^k Z_i^2 - \sum_{i=1}^d Z_i^2 = \sum_{i=d+1}^k Z_i^2 \sim \chi_{k-d}^2.$$

□

An important property, essential for all desirable tests, is that the probability of type II error of the test at any level  $\alpha$  tends to 0 (ie, the power tends to 1), as  $n \rightarrow \infty$  for all departures from the null hypothesis. This property of a test is called *Consistency*.

**Theorem 7.2.2.** *The LRT is consistent, ie, for  $c > 0$ ,*

$$\lim_{n \rightarrow \infty} P_{\theta_0}[-2 \log \Lambda_n < c] = 0 \text{ if } \theta_0 \notin \Theta_0,$$

under regularity conditions (Cramér and Wald).

*Proof.* Suppose that the true  $\theta_0 \notin \Theta_0$ . Then arguing as in the proof of Theorem 7.2.1

$$\begin{aligned} -2 \log \Lambda_n &= -2 \left[ \sum_{i=1}^n l(X_i, \tilde{\theta}_n) - \sum_{i=1}^n l(X_i, \theta_0) \right] \\ &\quad + 2 \left[ \sum_{i=1}^n l(X_i, \hat{\theta}_n) - \sum_{i=1}^n l(X_i, \theta_0) \right] \\ &= -2 \left[ \sum_{i=1}^n l(X_i, \tilde{\theta}_n) - \sum_{i=1}^n l(X_i, \theta_0) \right] \\ &\quad + \left\{ n^{-1/2} \sum_{i=1}^n \mathbf{i}(X_i, \theta_0) \right\}^T \mathbf{I}(\theta_0)^{-1} \left\{ n^{-1/2} \sum_{i=1}^n \mathbf{i}(X_i, \theta_0) \right\} + Y_n \\ &:= -2 \left[ \sum_{i=1}^n l(X_i, \tilde{\theta}_n) - \sum_{i=1}^n l(X_i, \theta_0) \right] + b_n^2 + Y_n, \end{aligned}$$

where  $Y_n = o_P(1)$ .

Since

- (i)  $P[X_n + Y_n + b_n^2 < c] \leq P[X_n + Y_n < c] \leq P[X_n < c + |Y_n|] \leq P[X_n < 2c] + P[|Y_n| \geq c]$ ,
- (ii)  $\lim_{n \rightarrow \infty} P[X_n < c] = 0$  for all  $c > 0$ , and
- (iii)  $Y_n = o_P(1)$

together imply  $\lim_{n \rightarrow \infty} P[X_n + Y_n + b_n^2 < c] = 0$  for all  $c > 0$ , it is enough to show that for all  $c > 0$ ,

$$\begin{aligned} & \lim_{n \rightarrow \infty} P_{\theta_0} \left[ - \sum_{i=1}^n l(X_i, \tilde{\theta}_n) + \sum_{i=1}^n l(X_i, \theta_0) < c \right] \\ &= \lim_{n \rightarrow \infty} P_{\theta_0} \left[ \inf_{\theta \in \Theta_0} n^{-1} \sum_{i=1}^n \log\{f(X_i, \theta_0)/f(X_i, \theta)\} < c/n \right] = 0 \end{aligned}$$

to prove the theorem.

Referring to the notations introduced in Section 7.1.1 and the proof of Theorem 7.1.4, define iid rv's

$$Z_i(\theta_0, S) = \inf_{\theta \in S} \log\{f(X_i, \theta_0)/f(X_i, \theta)\}, \quad i = 1, 2, \dots$$

for any set  $S \subset \mathbb{R}^k$  so that  $\theta_0 \notin S$ . Then under Wald's conditions, there exist a finite number of sets  $S_0, S_1, \dots, S_r$  such that

- (i)  $\Theta_0 \subset \bigcup_{j=0}^r S_j$  and (ii) for  $\theta_0 \notin \Theta_0$ ,  $I(\theta_0, S_j) = E_{\theta_0}[Z(\theta_0, S_j)] > 0$ ,  $j = 0, \dots, r$ .

We demonstrate the existence of such sets by the following steps

- (a) Choose a neighborhood  $N_{\infty, j^*}$  of  $\infty$  such that

$$E[Z(\theta_0, N_{\infty, j^*} \cap \Theta_0)] \geq E[Z(\theta_0, N_{\infty, j^*})] > 0.$$

- (b)  $N_{\infty, j^*}^c \cap \Theta_0$  is a closed bounded set in  $\mathbb{R}^k$  and therefore, has the Heine-Borel property of having a finite subcover for every open cover.
- (c) For every  $\theta \in N_{\infty, j^*}^c \cap \Theta_0$ , choose  $N_{\theta, j}$  such that  $E[Z(\theta_0, N_{\theta, j})] > 0$ . These sets form an open cover of  $N_{\infty, j^*}^c \cap \Theta_0$ , from which we now choose a finite subcover  $N_{\theta_1, j_1}, \dots, N_{\theta_r, j_r}$ .
- (d) Call these sets  $S_1 = N_{\theta_1, j_1}, \dots, S_r = N_{\theta_r, j_r}$  and let  $S_0 = N_{\infty, j^*} \cap \Theta_0$ . Then  $S_0, S_1, \dots, S_r$  cover  $\Theta_0$  and satisfy (i) and (ii) above. The choice of  $N_{\infty, j^*}$  and  $N_{\theta, j}$  for every  $\theta \in N_{\infty, j^*}^c \cap \Theta_0$  with these properties is ensured by Wald's regularity conditions.

It now follows that

$$\begin{aligned} & \inf_{\theta \in \Theta_0} n^{-1} \sum_{i=1}^n \log\{f(X_i, \theta_0)/f(X_i, \theta)\} \\ & \geq \min_{0 \leq j \leq r} \inf_{\theta \in S_j} n^{-1} \sum_{i=1}^n \log\{f(X_i, \theta_0)/f(X_i, \theta)\} \\ & \geq \min_{0 \leq j \leq r} n^{-1} \sum_{i=1}^n \inf_{\theta \in S_j} \log\{f(X_i, \theta_0)/f(X_i, \theta)\} = \min_{0 \leq j \leq r} n^{-1} \sum_{i=1}^n Z_i(\theta_0, S_j). \end{aligned}$$

Thus

$$\begin{aligned} P_{\theta_0} & \left[ \inf_{\theta \in \Theta_0} n^{-1} \sum_{i=1}^n \log\{f(X_i, \theta_0)/f(X_i, \theta)\} < c/n \right] \\ & = P_{\theta_0} \left[ \min_{0 \leq j \leq r} n^{-1} \sum_{i=1}^n Z_i(\theta_0, S_j) < c/n \right] \leq \sum_{j=0}^r P_{\theta_0} \left[ n^{-1} \sum_{i=1}^n Z_i(\theta_0, S_j) < c/n \right], \end{aligned}$$

which tends to 0 as  $n \rightarrow \infty$ , because

$$n^{-1} \sum_{i=1}^n Z_i(\theta_0, S_j) \xrightarrow{a.s.} E_{\theta_0}[Z(\theta_0, S_j)] = I(\theta_0, S_j) > 0.$$

□

**Example 7.2.1** (Homogeneity of Exponential Distributions). Let  $(X_{i1}, \dots, X_{in_i})$ ,  $i = 1, \dots, k$ , be independent random samples from  $\text{Exponential}(\theta_i)$ ,  $i = 1, \dots, k$ , respectively. Find the LRT for  $H_0: \theta_1 = \dots = \theta_k$  against all possible departures from  $H_0$ .

*Solution.* Let  $\hat{\theta}_{ni}$  and  $\bar{\theta}_N$  be, respectively, the unrestricted MLE and the restricted MLE under  $H_0$  of  $\theta_i$ . Then

$$\begin{aligned} \hat{\theta}_{ni} & = \bar{T}_i = T_i/n_i \text{ and } \bar{\theta}_N = \bar{T} = \sum_{i=1}^k T_i / \sum_{i=1}^k n_i, \text{ where} \\ T_i & = \sum_{j=1}^{n_i} X_{ij}. \end{aligned}$$

Now the LRT statistic is

$$\begin{aligned} \Lambda_N & = \prod_{i=1}^k \prod_{j=1}^{n_i} f(X_{ij}, \bar{\theta}_N) / \prod_{i=1}^k \prod_{j=1}^{n_i} f(X_{ij}, \hat{\theta}_{ni}) \\ & = \left[ \frac{1}{\bar{\theta}_N^N} e^{-T/\bar{\theta}_N} \right] / \left[ \prod_{i=1}^k \frac{1}{\hat{\theta}_{ni}^{n_i}} e^{-T_i/\hat{\theta}_{ni}} \right] = \prod_{i=1}^k \left( \frac{\hat{\theta}_{ni}}{\bar{\theta}_N} \right)^{n_i}. \end{aligned}$$

Since  $\hat{\theta}_{ni}/\bar{\theta}_N = 1 + o_P(1)$  under  $H_0$  and since  $\log(1 + Y_N) = Y_N - (1/2)Y_N^2\{1 + o_P(1)\}$  if  $Y_N = o_P(1)$ , it follows that

$$\begin{aligned} -2 \log \Lambda_N & = -2 \sum_{i=1}^k n_i \log(\hat{\theta}_{ni}/\bar{\theta}_N) \\ & = -2 \sum_{i=1}^k n_i [(\hat{\theta}_{ni}/\bar{\theta}_N - 1) - (1/2)(\hat{\theta}_{ni}/\bar{\theta}_N - 1)^2\{1 + o_P(1)\}] \end{aligned}$$

under  $H_0$ . It is easy to check that  $\sum_{i=1}^k n_i(\hat{\theta}_{ni}/\bar{\theta}_N - 1) = 0$ . Hence

$$-2 \log \Lambda_N = \sum_{i=1}^k \left[ n_i(\hat{\theta}_{ni} - \bar{\theta}_N)^2 / \bar{\theta}_N^2 \right] \{1 + o_P(1)\} \xrightarrow{\mathcal{L}} \chi_{k-1}^2$$

under  $H_0$  by [Theorem 7.2.1](#). Moreover, since  $\sqrt{n_i}(\hat{\theta}_{ni} - \theta)/\theta \xrightarrow{\mathcal{L}} N(0, 1)$  independently for each  $i$  and since  $\bar{\theta}_N = \theta + o_P(1)$ , we can argue more directly, that

$$\begin{aligned} & \sum_{i=1}^k \left[ n_i(\hat{\theta}_{ni} - \bar{\theta}_N)^2 / \bar{\theta}_N^2 \right] \{1 + o_P(1)\} \\ &= \sum_{i=1}^k n_i(\hat{\theta}_{ni} - \bar{\theta}_N)^2 / \bar{\theta}_N^2 + o_P(1) \xrightarrow{\mathcal{L}} \chi_{k-1}^2. \end{aligned}$$

**Example 7.2.2** (Homogeneity of Multinomial Probabilities). Consider  $m$  mutually exclusive and exhaustive categories and let

$$\pi_i = (\pi_{i1}, \dots, \pi_{im}), \quad i = 1, \dots, k, \quad \pi_{ij} > 0 \text{ and } \sum_{j=1}^m \pi_{ij} = 1$$

be the probability distributions over these categories in  $k$  populations. These probability vectors are called multinomial probabilities. Let  $(n_{i1}, \dots, n_{im})$  be the frequencies in the  $m$  categories in independent random samples of sizes  $n_i$ ,  $i = 1, \dots, k$ , from these populations. Find the LRT for  $H_0$ :  $\pi_1 = \dots = \pi_k$  against all possible departures from  $H_0$ .

*Solution.* The unrestricted MLEs and restricted MLEs under  $H_0$  of  $\pi_i$  are

$$\hat{\pi}_i = (\hat{\pi}_{i1}, \dots, \hat{\pi}_{im}) \text{ and } \hat{\pi}_i^0 = (\hat{\pi}_{i1}^0, \dots, \hat{\pi}_{im}^0),$$

respectively, where  $\hat{\pi}_{ij} = n_{ij}/n_i$  and  $\hat{\pi}_{ij}^0 = n_{0j}/N$ , with  $n_{0j} = \sum_{i=1}^k n_{ij}$  and  $N = \sum_{i=1}^k n_i$ . For notational convenience, we also let  $n_i = \sum_{j=1}^m n_{ij} = n_{i0}$ . The LRT statistic can now be written as

$$\begin{aligned} \Lambda_N &= \prod_{i=1}^k f((n_{i1}, \dots, n_{im}), \hat{\pi}_i^0) / \prod_{i=1}^k f((n_{i1}, \dots, n_{im}), \hat{\pi}_i) \\ &= \prod_{i=1}^k \prod_{j=1}^m (\hat{\pi}_{ij}^0 / \hat{\pi}_{ij})^{n_{ij}}, \end{aligned}$$

cancelling the multinomial coefficients from the numerator and denominator. Again, using  $\hat{\pi}_{ij}^0 / \hat{\pi}_{ij} = 1 + o_P(1)$  and the property of  $\log(1 + Y_n)$  for  $Y_n = o_P(1)$  as in [Example 7.1.1](#), we have

$$\begin{aligned}
 -2 \log \Lambda_N &= -2 \sum_{i=1}^k \sum_{j=1}^m n_{ij} \log(\hat{\pi}_{ij}^0 / \hat{\pi}_{ij}) \\
 &= -2 \sum_{i=1}^k \sum_{j=1}^m n_{i0} \hat{\pi}_{ij} \left[ \frac{\hat{\pi}_{ij}^0 - \hat{\pi}_{ij}}{\hat{\pi}_{ij}} - (1/2) \left( \frac{\hat{\pi}_{ij}^0 - \hat{\pi}_{ij}}{\hat{\pi}_{ij}} \right)^2 \{1 + o_P(1)\} \right] \\
 &= -2 \sum_{i=1}^k \sum_{j=1}^m n_{i0} (\hat{\pi}_{ij}^0 - \hat{\pi}_{ij}) + \sum_{i=1}^k \sum_{j=1}^m n_{i0} \frac{(\hat{\pi}_{ij}^0 - \hat{\pi}_{ij})^2}{\hat{\pi}_{ij}} \{1 + o_P(1)\}
 \end{aligned}$$

under  $H_0$ . Since  $\sum_{i=1}^k n_{i0} (\hat{\pi}_{ij}^0 - \hat{\pi}_{ij}) = 0$  for each  $j$ , and

$$\begin{aligned}
 \sum_{i=1}^k \sum_{j=1}^m n_{i0} \frac{(\hat{\pi}_{ij}^0 - \hat{\pi}_{ij})^2}{\hat{\pi}_{ij}} &= \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - (n_{i0} n_{0j}/N))^2}{n_{ij}} \\
 &= \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - (n_{i0} n_{0j}/N))^2}{(n_{i0} n_{0j}/N)} \{1 + o_P(1)\}
 \end{aligned}$$

under  $H_0$ , we finally have

$$-2 \log \Lambda_N = \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - (n_{i0} n_{0j}/N))^2}{(n_{i0} n_{0j}/N)} + o_P(1) \xrightarrow{\mathcal{L}} \chi_{(k-1)(m-1)}^2.$$

### 7.3 Asymptotic Properties of MLE and LRT Based on Independent Nonidentically Distributed Data

In this section, we shall briefly indicate the behavior of MLE and LRT when the observations are independent but nonidentically distributed.

Suppose that  $(X_{11}, \dots, X_{1n_1}), (X_{21}, \dots, X_{2n_2}), \dots, (X_{m1}, \dots, X_{mn_m})$  are independent random samples from distributions with pdf/pmf  $f_1(x, \theta), \dots, f_m(x, \theta)$ , respectively, and let  $n = n_1 + \dots + n_m$  be the total sample size.

We assume that the regularity conditions stated earlier, hold for each  $\{f_j(x, \theta), \theta \in \Theta \subset \mathbb{R}^k\}$  and let

$$I_{j,rs}(\theta_0) = E[-\ddot{I}_{j,rs}(X, \theta_0)], \quad I_j(\theta_0) = ((I_{j,rs}(\theta_0))), \quad j = 1, \dots, m.$$

Let  $\hat{\theta}_n$  be the MLE of  $\theta$  based on the pooled data consisting of all  $X_{ji}$ ,  $i = 1, \dots, n_j$ ,  $j = 1, \dots, m$ .

If  $n_1, \dots, n_m \rightarrow \infty$  in such a way that  $n_j/n \rightarrow \lambda_j > 0$  for all  $j$  (of course  $\lambda_1 + \dots + \lambda_m = 1$ ), and if  $\mathbf{I}(\theta_0) = \sum_{j=1}^m \lambda_j \mathbf{I}_j(\theta_0)$  is positive definite (even if some of the individual  $\mathbf{I}_j(\theta_0)$  is singular), then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{\mathcal{L}} \mathbf{W} \sim N_k(\mathbf{0}, \mathbf{I}(\theta_0)^{-1}).$$

Asymptotic properties of the LRT discussed earlier for the iid case, also hold in the non-iid case described above.

## 7.4 Frequency $\chi^2$

An important class of hypothesis testing problems arises in the context of the multinomial distribution  $Multi(n; \pi_1, \dots, \pi_m)$ , where  $\pi_j \geq 0$  with  $\sum \pi_j = 1$  are the probabilities of a random sample belonging to the classes  $(1, \dots, m)$ . The data consist of a random sample  $(X_1, \dots, X_n)$  where the  $X_i$ 's are  $m$ -dim iid rv's with

$$f(\mathbf{e}_j, \boldsymbol{\pi}) = P(X_i = \mathbf{e}_j) = \pi_j, \quad j = 1, \dots, m,$$

where  $\mathbf{e}_j$  is the  $m$ -dim vector with 1 at the  $j$ th coordinate and 0 for the rest of the coordinates, and  $\boldsymbol{\pi}^T = (\pi_1, \dots, \pi_m)$ . As in [Section 7.1.3](#), let

$$(n_1, \dots, n_m)^T = \sum_{i=1}^n \mathbf{X}_i,$$

$$n_j = \sum_{i=1}^n I(X_i = \mathbf{e}_j) = \text{observed frequency in the } j\text{th class.}$$

For testing  $H_0: \pi_j = \pi_{j0}, j = 1, \dots, m$ , with  $\pi_{j0} \geq 0$  with  $\sum_{j=1}^m \pi_{j0} = 1$ , against all possible alternatives satisfying these constraints, the test statistic

$$T_{0n} = \sum_{j=1}^m \frac{(n_j - n\pi_{j0})^2}{n\pi_{j0}}, \tag{7a}$$

known as the frequency  $\chi^2$ , introduced by K. Pearson, is widely used. In many problems of practical importance such as testing for independence in  $r \times s$  contingency tables, the null hypothesis value for the  $j$ th class is  $\pi_{j0}(\theta_1, \dots, \theta_k)$  where  $\pi_{10}(\cdot), \dots, \pi_{m0}(\cdot)$  are known functions of an unknown  $k$ -dim parameter  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^T$ . The test statistic used to test such a composite hypothesis  $H_1: \pi_j = \pi_{j0}(\boldsymbol{\theta}), j = 1, \dots, m$  for given functions  $\boldsymbol{\pi}_0(\cdot) = (\pi_{10}(\cdot), \dots, \pi_{m0}(\cdot))^T$  of an unknown  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^k$  against all possible alternatives is

$$T_{1n} = \sum_{j=1}^m \frac{[n_j - n\pi_{j0}(\hat{\boldsymbol{\theta}}_n)]^2}{n\pi_{j0}(\hat{\boldsymbol{\theta}}_n)}, \tag{7b}$$

where  $\hat{\boldsymbol{\theta}}_n = (\hat{\theta}_{n1}, \dots, \hat{\theta}_{nk})^T$  is the MLE of  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^T$ . In both the above problems, we reject the respective null hypotheses if  $T_{0n}$  or  $T_{1n}$  is too large. Our goal is to find the asymptotic distribution of  $T_{0n}$  under  $H_0$  and of  $T_{1n}$  under  $H_1$  so that the critical values  $c_{0\alpha}$  and  $c_{1\alpha}$  for level  $\alpha$  tests can be determined so that

$$P_{H_0}[T_{0n} \geq c_{0\alpha}] \approx \alpha \text{ and } P_{H_1}[T_{1n} \geq c_{1\alpha}] \approx \alpha,$$

for large  $n$ . This will be accomplished in two steps:

1. Simplifying  $T_{0n}$  given by Eq. (7a) and  $T_{1n}$  given by Eq. (7b) to forms more convenient for asymptotics.
2. Deriving the asymptotic distributions of these simplified forms of  $T_{0n}$  and  $T_{1n}$ .

The statistic  $T_{0n}$  is called the frequency  $\chi^2$  because its null distribution follows the  $\chi^2$  distribution with  $m - 1$  df and the statistic  $T_{1n}$  also has a  $\chi^2$  distribution but with  $m - k - 1$  df. These asymptotic distributions will be derived in this section.

We first define some vectors and matrices which will be useful in handling  $T_{0n}$  and  $T_{1n}$ . In dealing with  $T_{0n}$ , let

$$q_j = \sqrt{\pi_{j0}}, \mathbf{q}^T = (q_1, \dots, q_m), \boldsymbol{\Lambda}_q = \text{diag}((q_j)), \boldsymbol{\Lambda}_{1/q} = \boldsymbol{\Lambda}_q^{-1} = \text{diag}((1/q_j)).$$

Then  $\text{diag}((\pi_{j0})) = \boldsymbol{\Lambda}_q^2$ ,  $\boldsymbol{\pi}_0 = \boldsymbol{\Lambda}_q \mathbf{q}$ , and under  $H_0$ ,

$$\mathbb{E}[\mathbf{X}] = \boldsymbol{\pi}_0 = \boldsymbol{\Lambda}_q \mathbf{q} \text{ and } \boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{X} - \boldsymbol{\pi}_0)(\mathbf{X} - \boldsymbol{\pi}_0)^T] = \boldsymbol{\Lambda}_q(\mathbf{I} - \mathbf{q}\mathbf{q}^T)\boldsymbol{\Lambda}_q \quad (8)$$

after some simplification. Next let

$$\mathbf{Z}_n = n^{-1/2}(n_1 - n\pi_{10}, \dots, n_m - n\pi_{m0})^T = n^{-1/2} \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\pi}_0).$$

Then

$$\mathbf{Z}_n = n^{-1/2} \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}(\mathbf{X}_i)) \xrightarrow{\mathcal{L}} \mathbf{Z} \sim N_m(\mathbf{0}, \boldsymbol{\Sigma}) \text{ under } H_0. \quad (9)$$

**Lemma 7.4.1.** Under  $H_0$ ,  $T_{0n} \xrightarrow{\mathcal{L}} \mathbf{U}^T \mathbf{U}$ , where  $\mathbf{U} \sim N_m(\mathbf{0}, \mathbf{I} - \mathbf{q}\mathbf{q}^T)$ .

*Proof.* Using Eq. (8)  $T_{0n}$  can be written as

$$\begin{aligned} T_{0n} &= \sum_{j=1}^m \left\{ (\sqrt{n}q_j)^{-1} (n_j - n\pi_{j0}) \right\}^2 = \sum_{j=1}^m (1/q_j^2) \left\{ n^{-1/2} (n_j - n\pi_{j0}) \right\}^2 \\ &= (\boldsymbol{\Lambda}_{1/q} \mathbf{Z}_n)^T (\boldsymbol{\Lambda}_{1/q} \mathbf{Z}_n) := \mathbf{U}_n^T \mathbf{U}_n \xrightarrow{\mathcal{L}} \mathbf{U}^T \mathbf{U}, \end{aligned}$$

where  $\mathbf{U}_n = \boldsymbol{\Lambda}_{1/q} \mathbf{Z}_n \xrightarrow{\mathcal{L}} \boldsymbol{\Lambda}_{1/q} \mathbf{Z} = \mathbf{U} \sim N_m(\mathbf{0}, \mathbf{I} - \mathbf{q}\mathbf{q}^T)$ , because the covariance matrix of  $\mathbf{U}$  is  $\boldsymbol{\Lambda}_{1/q} \boldsymbol{\Sigma} \boldsymbol{\Lambda}_{1/q} = \mathbf{I} - \mathbf{q}\mathbf{q}^T$  after some simplification.  $\square$

In dealing with  $T_{1n}$ , we continue to use the above notations with  $\pi_{j0}(\theta_0)$  instead of  $\pi_{j0}$ . The problem with  $T_{1n}$  is the presence of  $\hat{\theta}_n$  in both the numerator and the denominator. Replacing  $\hat{\theta}_n$  by  $\theta_0$  in the denominator is relatively easy and will be taken care at first. Handling  $\hat{\theta}_n$  in the numerator will take more work. We start with  $T_{1n}$  and after algebraic re-arrangements, write

$$\begin{aligned}
T_{1n} &= \sum_{j=1}^m [n_j - n\pi_{j0}(\hat{\theta}_n)]^2 / (n\pi_{j0}(\hat{\theta}_n)) \\
&= \sum_{j=1}^m [n_j - n\pi_{j0}(\hat{\theta}_n)]^2 / (n\pi_{j0}(\theta_0)) \{1 - R_{nj}\} := T_{1n}^* + R_n, \text{ where} \\
T_{1n}^* &= \sum_{j=1}^m [n_j - n\pi_{j0}(\hat{\theta}_n)]^2 / (n\pi_{j0}(\theta_0)) \text{ and} \\
|R_n| &\leq T_{1n}^* \max_{1 \leq j \leq m} |R_{nj}| = T_{1n}^* o_P(1),
\end{aligned}$$

because under regularity conditions,  $\pi_{j0}(\hat{\theta}_n) - \pi_{j0}(\theta_0) = o_P(1)$ . We shall now show that  $T_{1n}^* \xrightarrow{\mathcal{L}} T_1^*$  (to be determined), so that  $T_{1n}^* = O_P(1)$  which would imply  $R_n = O_P(1)o_P(1) = o_P(1)$ , proving that  $T_{1n} = T_{1n}^* + o_P(1) \xrightarrow{\mathcal{L}} T_1^*$ .

We now work on

$$\begin{aligned}
T_{1n}^* &= \sum_{j=1}^m [n_j - n\pi_{j0}(\hat{\theta}_n)]^2 / (n\pi_{j0}(\theta_0)) \\
&= \sum_{j=1}^m \left[ \frac{n_j - n\pi_{j0}(\theta_0)}{\sqrt{n\pi_{j0}(\theta_0)}} - \sqrt{n} \frac{\pi_{j0}(\hat{\theta}_n) - \pi_{j0}(\theta_0)}{\sqrt{\pi_{j0}(\theta_0)}} \right]^2 \\
&= (\mathbf{U}_n - \mathbf{W}_n)^T (\mathbf{U}_n - \mathbf{W}_n), \tag{10}
\end{aligned}$$

where  $\mathbf{U}_n = \Lambda_{1/q} \mathbf{Z}_n$  as in Eq. (9) and the proof of Lemma 7.4.1 with  $\pi_0(\theta_0)$  in place of  $\pi_0$ , and

$$\mathbf{W}_n^T = \sqrt{n} \left( \frac{\pi_{10}(\hat{\theta}_n) - \pi_{10}(\theta_0)}{\sqrt{\pi_{10}(\theta_0)}}, \dots, \frac{\pi_{m0}(\hat{\theta}_n) - \pi_{m0}(\theta_0)}{\sqrt{\pi_{m0}(\theta_0)}} \right).$$

Assuming regularity conditions including existence of  $\partial\pi_{j0}/\partial\theta_r = \dot{\pi}_{rj}$  for  $1 \leq r \leq k$ ,  $1 \leq j \leq m$  and positive-definiteness of  $\mathbf{I}(\theta_0) = ((I_{rs}(\theta_0)))$ , let  $d_{rj} = \dot{\pi}_{rj}/q_j$ ,  $\mathbf{D} = ((d_{rj}))_{k \times m}$  and  $I_{rs}(\theta_0) = \sum_{j=1}^m d_{rj} d_{sj} = (r, s)$ th element of  $\mathbf{D}\mathbf{D}^T$ .

In these notations,

$$\pi_{j0}(\hat{\theta}_n) - \pi_{j0}(\theta_0) = \sum_{r=1}^k \dot{\pi}_{rj} (\hat{\theta}_{nr} - \theta_{0r}) + o_P(1), \quad 1 \leq j \leq m,$$

and now  $\mathbf{W}_n^T$  can be expressed as

$$\begin{aligned}
\mathbf{W}_n^T &= \sqrt{n}(\hat{\theta}_n - \theta_0)^T \mathbf{D} + o_P(1) = \left\{ \mathbf{I}(\theta_0)^{-1} n^{-1/2} \sum_{i=1}^n \dot{\mathbf{l}}(\mathbf{X}_i, \theta_0) \right\}^T \mathbf{D} + o_P(1) \\
&= \{I(\theta_0)^{-1} \mathbf{D} \mathbf{U}_n\}^T \mathbf{D} + o_P(1), \tag{11}
\end{aligned}$$

using the expression for  $\sqrt{n}(\hat{\theta}_n - \theta_0)$  obtained in the course of the proof of [Theorem 7.1.3](#), because

$$\begin{aligned}\sum_{i=1}^n \dot{l}(X_i, \theta_0) &= \sum_{i=1}^n \sum_{j=1}^m I[X_i = e_j] \dot{l}(e_j, \theta_0) = \sum_{j=1}^m n_j \dot{l}(e_j, \theta_0) \\ &= \sum_{j=1}^m \frac{n_j}{\sqrt{\pi_j(\theta_0)}} \left( \frac{\dot{\pi}_{1j}(\theta_0)}{q_j}, \dots, \frac{\dot{\pi}_{kj}(\theta_0)}{q_j} \right)^T \\ &= \sqrt{n} \sum_{j=1}^m \frac{n_j - n\pi_{j0}(\theta_0)}{\sqrt{n\pi_j(\theta_0)}} \left( \frac{\dot{\pi}_{1j}(\theta_0)}{q_j}, \dots, \frac{\dot{\pi}_{kj}(\theta_0)}{q_j} \right)^T \\ &= \sqrt{n} \sum_{j=1}^m U_{nj}(d_{1j}, \dots, d_{kj}) = \sqrt{n} \mathbf{D}\mathbf{U}_n,\end{aligned}$$

since  $\sum_{j=1}^m \pi_{j0}(\theta_0) = 1$  implies  $\sum_{j=1}^m \dot{\pi}_{rj}(\theta_0) = 0$ ,  $r = 1, \dots, k$ .

From Eq. (11), it now follows that under  $H_1$ ,

$$\begin{aligned}\mathbf{V}_n &= \mathbf{U}_n - \mathbf{W}_n = \mathbf{U}_n - \mathbf{D}^T \{\mathbf{I}(\theta_0)^{-1} \mathbf{D}\mathbf{U}_n\} = [\mathbf{I} - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D}] \mathbf{U}_n \\ &\xrightarrow{\mathcal{L}} \mathbf{V} = [\mathbf{I} - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D}] \mathbf{U} \sim N_m(\mathbf{0}, \Sigma_V), \text{ where} \\ \Sigma_V &= [\mathbf{I} - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D}] (\mathbf{I} - \mathbf{q}\mathbf{q}^T) [\mathbf{I} - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D}] \\ &= \mathbf{I} - \mathbf{q}\mathbf{q}^T - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D}.\end{aligned}$$

**Lemma 7.4.2.** Under  $H_1$ ,  $T_{1n} \xrightarrow{\mathcal{L}} \mathbf{V}^T \mathbf{V}$  where  $\mathbf{V} \sim N_m(\mathbf{0}, \mathbf{I} - \mathbf{q}\mathbf{q}^T - \mathbf{D}^T (\mathbf{D}\mathbf{D}^T)^{-1} \mathbf{D})$ .

*Proof.* As shown in Eq. (10)

$$\begin{aligned}T_{1n}^* &= (\mathbf{U}_n - \mathbf{W}_n)^T (\mathbf{U}_n - \mathbf{W}_n) = \mathbf{V}_n^T \mathbf{V}_n \\ &\xrightarrow{\mathcal{L}} T_1^* = \mathbf{V}^T \mathbf{V}.\end{aligned}$$

Hence

$$T_{1n} = T_{1n}^* + R_n \xrightarrow{\mathcal{L}} T_1^* = \mathbf{V}^T \mathbf{V},$$

because  $|R_n| \leq T_n^* \max_{1 \leq j \leq m} |R_{nj}| = o_p(1)$  as observed earlier.  $\square$

The distributions of  $\mathbf{U}^T \mathbf{U}$  and  $\mathbf{V}^T \mathbf{V}$  are obtained by using the following two lemmas.

**Lemma 7.4.3.** If  $\mathbf{Y} \sim N_m(\mathbf{0}, \mathbf{C})$ , then  $\mathbf{Y}^T \mathbf{Y} \stackrel{\mathcal{D}}{\equiv} \sum_{j=1}^m \lambda_j \xi_j^2$  where  $\lambda_1, \dots, \lambda_m$  are the eigenvalues of  $\mathbf{C}$  and  $\xi_1, \dots, \xi_m$  are iid  $N(0, 1)$ .

*Proof.* Let  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]$  and  $\Lambda = \text{diag}((\lambda_j))$ , where  $\mathbf{a}_1, \dots, \mathbf{a}_m$  are orthonormal eigenvectors and  $\lambda_1, \dots, \lambda_m$  the corresponding eigenvalues of  $\mathbf{C}$ . Then  $\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$  and  $\mathbf{C} = \mathbf{A} \Lambda \mathbf{A}^T$ . Let  $\mathbf{W} = \mathbf{A}^T \mathbf{Y}$ . Then  $\mathbf{W}^T \mathbf{W} = \mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{W} \sim N_m(\mathbf{0}, \mathbf{A}^T (\mathbf{A} \Lambda \mathbf{A}^T) \mathbf{A}) = N_m(\mathbf{0}, \Lambda)$ , so that

$$\mathbf{Y}^T \mathbf{Y} = \mathbf{W}^T \mathbf{W} = \sum_{j=1}^m W_j^2 = \sum_{j=1}^m \lambda_j \left( W_j / \sqrt{\lambda_j} \right)^2 = \sum_{j=1}^m \lambda_j \xi_j^2,$$

where  $\xi_1, \dots, \xi_m$  are iid  $N(0, 1)$ .  $\square$

**Lemma 7.4.4.** (i) The matrix  $\mathbf{I} - \mathbf{q}\mathbf{q}^T$  has  $m-1$  eigenvalues equal to 1 and one eigenvalue equal to 0. (ii) The matrix  $\mathbf{I} - \mathbf{q}\mathbf{q}^T - \mathbf{D}^T(\mathbf{D}\mathbf{D}^T)^{-1}\mathbf{D}$  has  $m-k-1$  eigenvalues equal to 1 and  $k+1$  eigenvalues equal to 0.

*Proof.* First note that  $\lambda_1, \dots, \lambda_m$  are the eigenvalues of  $\mathbf{C}$  iff  $1 - \lambda_1, \dots, 1 - \lambda_m$  are the eigenvalues of  $\mathbf{I} - \mathbf{C}$ , because  $(\mathbf{I} - \mathbf{C})\mathbf{a} = (1 - \lambda)\mathbf{a} \iff \mathbf{C}\mathbf{a} = \lambda\mathbf{a}$ . Therefore, we need to consider eigenvalues of  $\mathbf{q}\mathbf{q}^T$  and  $\mathbf{q}\mathbf{q}^T + \mathbf{D}^T(\mathbf{D}\mathbf{D}^T)^{-1}\mathbf{D}$ . First,  $(\mathbf{q}\mathbf{q}^T)\mathbf{q} = \mathbf{q}(\mathbf{q}^T\mathbf{q}) = \mathbf{q}$  (since  $\mathbf{q}^T\mathbf{q} = \sum_{j=1}^m \pi_j 0 = 1$ ), so 1 is an eigenvalue of  $\mathbf{q}\mathbf{q}^T$  with eigenvector  $\mathbf{q}$ . Now let  $\mathbf{a}_1, \dots, \mathbf{a}_{m-1}$  be the other  $m-1$  eigenvectors of  $\mathbf{q}\mathbf{q}^T$ . Then  $\mathbf{q}^T\mathbf{a}_i = 0$ ,  $i = 1, \dots, m-1$ , so  $(\mathbf{q}\mathbf{q}^T)\mathbf{a}_i = \mathbf{q}(\mathbf{q}^T\mathbf{a}_i) = 0$  for all  $i$ , showing that the other  $m-1$  eigenvalues of  $\mathbf{q}\mathbf{q}^T$  are 0. This proves (i). Similarly, we can verify that  $\mathbf{q}$  and the  $k$  row vectors of  $\mathbf{D}$  are the eigenvectors of  $\mathbf{q}\mathbf{q}^T + \mathbf{D}^T(\mathbf{D}\mathbf{D}^T)^{-1}\mathbf{D}$  with corresponding eigenvalues all equal to 1 and other eigenvalues all equal to 0. This proves (ii).  $\square$

**Theorem 7.4.1.** Under  $H_0$ ,  $T_{0n} \xrightarrow{\mathcal{L}} \mathbf{U}^T \mathbf{U} \sim \chi_{m-1}^2$  and under  $H_1$ ,  $T_{1n} \xrightarrow{\mathcal{L}} \mathbf{V}^T \mathbf{V} \sim \chi_{m-k-1}^2$ .

*Proof.* Use Lemmas 7.4.1 and 7.4.4.  $\square$

**Example 7.4.1.** The number of flight cancelations between 6 am and noon at a certain airport was recorded for each day between the months of March and June. The data are summarized below with  $k$  = number of cancelations and  $n_k$  = number of days with  $k$  cancelations between 6 am and noon.

$k$	0	1	2	3	4	Total
$n_k$	35	55	22	8	2	$n = 122$

Test at level  $\alpha = 0.05$  whether the number of flight cancelations follows a Poisson distribution.

*Solution.* Under the null hypothesis  $H_0$  that the number of cancelations follows a Poisson distribution with an unspecified mean  $\theta$ , the MLE of  $\theta$  is

$$\hat{\theta}_n = \sum_k kn_k/n = 1.074 \text{ and } e^{-\hat{\theta}_n} = 0.3417.$$

The expected frequencies under  $H_0$  with  $Poisson(\hat{\theta}_n)$  are given below, together with the observed frequencies. Since the frequency  $\chi^2$  test is a large sample test and in the data, the frequency  $n_k = 2$  for  $k = 4$  is too small, we have pooled the classes  $k = 3$  and  $k = 4$ . (As a rule of thumb, we need the expected frequency in each class to be at least 5.)

$k$	0	1	2	$\geq 3$	Total
Observed frequency: $n_k$	35	55	22	10	122
Expected frequency: $nf(k, \hat{\theta}_n)$	41.69	44.77	24.04	11.52	122

From the observed and expected frequencies, the test statistic is obtained as

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}} = 1.07 + 2.34 + 0.18 + 0.19 = 3.78.$$

Under  $H_0$ , the test statistic is (asymptotically) distributed as  $\chi^2$  with

$$\text{df} = \text{number of classes} - \text{number of parameters estimated} - 1 = 4 - 1 - 1 = 2.$$

From the table of  $\chi^2$  distribution we have the critical value  $\chi_{2,0.05}^2 = 5.99$  (ie,  $P[\chi_2^2 > 5.99] = 0.05$ ). Since the observed value of  $\chi^2$  is smaller than the critical value, we accept  $H_0$  at level of significance  $\alpha = 0.05$ . In other words, at level  $\alpha = 0.05$ , the data indicate that the Poisson model is an acceptable fit.

*Remark 7.4.1.* Strictly speaking, having pooled some classes, one should calculate the MLE of the unspecified parameter(s) from the data after pooling. Since the calculation of MLE from pooled data is complicated, this issue is overlooked in practice. However, the use of MLE of  $\theta$  from original data in calculating the Frequency  $\chi^2$  statistic from pooled data results in the asymptotic distribution of the test statistic to be stochastically larger than a  $\chi^2$  with prescribed distribution. This error is not serious in fitting a Poisson distribution, but may be so in fitting a normal distribution based on frequencies in class intervals and using MLEs of  $\mu$  and  $\sigma^2$  from raw data (see [32]).

**Example 7.4.2** (Test for Independence in a Contingency Table). Let  $(A_1, \dots, A_k)$  and  $(B_1, \dots, B_m)$  be two classifications (into mutually exclusive and exhaustive categories) of a population with  $P(A_i) = \pi_{i0} > 0$ ,  $\sum_{i=1}^k \pi_{i0} = 1$  and  $P(B_j) = \pi_{0j} > 0$ ,  $\sum_{j=1}^m \pi_{0j} = 1$ . Also let  $P(A_i \cap B_j) = \pi_{ij}$ . In a random sample of  $N$  observations from this population, the frequency distribution over such a cross-classification is called a contingency table having the following layout in which the frequency of  $A_i B_j$  is  $n_{ij}$ .

Contingency Table

	$B_1$		$B_j$		$B_m$	Subtotal
$A_1$	$n_{11}$	$\cdots$	$n_{1j}$	$\cdots$	$n_{1m}$	$n_{10}$
	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$A_i$	$n_{i1}$	$\cdots$	$n_{ij}$	$\cdots$	$n_{im}$	$n_{i0}$
	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$A_k$	$n_{k1}$	$\cdots$	$n_{kj}$	$\cdots$	$n_{km}$	$n_{k0}$
Subtotal	$n_{01}$	$\cdots$	$n_{0j}$	$\cdots$	$n_{0m}$	$N$

In this model, we need to test  $H_0: \pi_{ij} = \pi_{i0}\pi_{0j}$  for all  $(i, j)$  (ie, the classifications  $\{A_i\}$  and  $\{B_j\}$  are independent), against all possible departures from  $H_0$ .

*Solution.* In this cross-classification, there are altogether  $km$  classes with observed frequencies  $n_{ij}$  and expected frequencies  $N\pi_{i0}\pi_{0j}$  under  $H_0$ . Since  $H_0$  involves unknown parameters  $\pi_{i0}$ ,  $i = 1, \dots, k$ , and  $\pi_{0j}$ ,  $j = 1, \dots, m$ , we use their MLEs  $\hat{\pi}_{i0} = n_{i0}/N$  and  $\hat{\pi}_{0j} = n_{0j}/N$ . However, due to the constraints  $\sum_{i=1}^k \pi_{i0} = 1 = \sum_{j=1}^m \pi_{0j}$ , only  $(k-1) + (m-1)$  parameters are actually estimated to find the expected frequencies  $N\hat{\pi}_{i0}\hat{\pi}_{0j} = n_{i0}n_{0j}/N$ . Thus the frequency  $\chi^2$  test statistic is

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{[n_{ij} - n_{i0}n_{0j}/N]^2}{n_{i0}n_{0j}/N},$$

which is asymptotically distributed under  $H_0$  as a  $\chi^2$  with

$$\begin{aligned} \text{df} &= \text{number of classes} - \text{number of parameters estimated} - 1 \\ &= km - \{(k-1) + (m-1)\} - 1 = (k-1)(m-1). \end{aligned}$$

Note that conditionally, given  $(n_{10}, \dots, n_{k0})$ , the problem is the same as the one discussed in [Example 7.2.2](#), where the LRT statistic is the same as the frequency  $\chi^2$  test statistics here, with the same asymptotic null distribution.

## Exercises

- 7.1. Let  $X_1, \dots, X_n$  be iid with pdf  $f(x, \theta) = \theta x^{-\theta-1} I_{(1, \infty)}(x)$ ,  $\theta > 2$ .
    - (a) Find the method of moments estimator  $\tilde{\theta}_n$  of  $\theta$ .
    - (b) Find the MLE  $\hat{\theta}_n$  of  $\theta$ .
    - (c) Find the asymptotic distributions of  $\sqrt{n}(\tilde{\theta}_n - \theta)$  and  $\sqrt{n}(\hat{\theta}_n - \theta)$ .
  - 7.2. Repeat Exercise 7.1 for the pdf  $f(x, \theta) = \theta(\theta+1)x^{\theta-1}(1-x)I_{(0,1)}(x)$ ,  $\theta > 0$ .
  - 7.3. Let  $X_1, \dots, X_n$  be iid with pdf
- $$f(x, \theta) = \theta I_{[0,1/3]}(x) + 2\theta I_{[1/3,2/3]}(x) + 3(1-\theta)I_{[2/3,1]}(x), \quad 0 < \theta < 1.$$
- (a) Find the MLE  $\hat{\theta}_n$  of  $\theta$  and its asymptotic distribution. Is  $\hat{\theta}_n$  unbiased?
  - (b) Is there a UMVUE of  $\theta$ ? If so, find it.
- 7.4. Let  $X_1, \dots, X_n$  be iid with pdf  $f(x, \theta) = \exp[-(x-\theta)]I_{[\theta, \infty)}(x)$ . Find the MLE  $\hat{\theta}_n$  of  $\theta$  and the asymptotic distribution of  $\hat{\theta}_n$  after appropriate normalization.
  - 7.5. Let  $(X_1, \dots, X_n)$  be a random sample from a log normal distribution with pdf

$$f(x, \theta) = \frac{1}{\sqrt{2\pi}x} \exp[-(\log x - \theta)^2/2]I_{(0, \infty)}(x).$$

This means,  $Y_i = \log X_i$ ,  $i = 1, \dots, n$ , are iid  $N(\theta, 1)$  so that

$(X_1, \dots, X_n) \stackrel{\mathcal{D}}{=} (e^{Z_1+\theta}, \dots, e^{Z_n+\theta})$ , where  $Z_1, \dots, Z_n$  are iid  $N(0, 1)$ .

- (a) Find the method of moments estimator  $\tilde{\theta}_n$  and the MLE  $\hat{\theta}_n$  of  $\theta$  based on  $X_1, \dots, X_n$ .

(b) Find the asymptotic distributions of  $\sqrt{n}(\tilde{\theta}_n - \theta)$  and  $\sqrt{n}(\hat{\theta}_n - \theta)$ .

(c) What is the asymptotic efficiency of  $\tilde{\theta}_n$ ?

- 7.6. Suppose that we observe  $U_{ij} = I_{(-\infty, a_i]}(X_{ij})$ ,  $j = 1, \dots, n_i$  and  $i = 1, \dots, k$ , where  $a_1 < \dots < a_k$  are known and  $X_{ij}$  are independent  $N(\mu, \sigma^2)$ . Let  $\boldsymbol{\theta}^\top = (\mu, \sigma^2)$ .

(a) Let  $\hat{\theta}_n$  be the MLE  $\boldsymbol{\theta}$  based on the observations  $\{U_{ij}\}$ . Find the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_n - \boldsymbol{\theta})$  where  $n = n_1 + \dots + n_k \rightarrow \infty$  in such a way that  $n_i/n \rightarrow \lambda_i \in (0, 1)$ .

(b) How would you calculate  $\hat{\theta}_n$  or an estimator asymptotically equivalent to  $\hat{\theta}_n$ ?  
[Let  $\hat{p}_i = \sum_{j=1}^{n_i} U_{ij}/n_i$ . Then  $\sqrt{n}(\hat{p}_i - \pi_i) = O_P(1)$ , where  $\pi_i = \Phi((a_i - \mu)/\sigma)$ .]

- 7.7. A discrete distribution with pmf

$$f(x, \theta) = \frac{\exp(-\theta)}{1 - \exp(-\theta)} \frac{\theta^x}{x!}, \quad x = 1, 2, \dots$$

is called a truncated Poisson distribution with parameter  $\theta$ . Consider  $\tilde{\theta}_n = 2v_2/v_1$ , where  $v_1$  and  $v_2$  are observed frequencies of  $X = 1$  and  $X = 2$  in a random sample of  $n$  observations from  $f(x, \theta)$  as an initial estimator  $\theta$  and construct a BAN estimator  $\hat{\theta}_n$  of  $\theta$ . Find the asymptotic distribution of  $\hat{\theta}_n$ .

- 7.8. Consider the multinomial distribution  $Multi(4, \boldsymbol{\pi}(\theta))$ , where  $\boldsymbol{\theta} = (p, q)$ ,  $\pi_O(\theta) = r^2$ ,  $\pi_A(\theta) = p^2 + 2pr$ ,  $\pi_B(\theta) = q^2 + 2qr$ ,  $\pi_{AB}(\theta) = 2pq$ ,  $p > 0$ ,  $q > 0$ , and  $r = 1 - p - q > 0$ . Suppose that in a random sample of size  $n$ , the cell frequencies are  $n_O$ ,  $n_A$ ,  $n_B$ , and  $n_{AB}$ . Set up the formulas for computing a BAN estimator of  $\boldsymbol{\theta}$ .

- 7.9. Let  $X_1, \dots, X_n$  be iid with pdf

$$f(x, \theta) = (1 - e^{-c/\theta})^{-1} (1/\theta) e^{-x/\theta} I_{(0, c]}(x), \quad \text{with a known } c,$$

and let  $\hat{\theta}_n$  be the MLE of  $\theta$  based on  $X_1, \dots, X_n$ . Since the likelihood equation does not have a closed form solution here, we try to obtain a BAN estimator, starting with an initial estimator  $\tilde{\theta}_{n0}$ . For this let

$$\tilde{\theta}_{n0} = \frac{c}{2 \log(p_n/q_n)}, \quad \text{where } p_n = n^{-1} \sum_{i=1}^n I_{(0, c/2]}(X_i) \text{ and } q_n = 1 - p_n.$$

Show that  $\tilde{\theta}_{n0} = \theta + op(1)$  and find the BAN estimator using  $\tilde{\theta}_{n0}$ . [Hint:

Let  $\phi = e^{-c/(2\theta)}$ . Then  $p_n$  is a sample proportion, estimating  
 $p = (1 - \phi)/(1 - \phi^2) = 1/(1 + \phi)$ .]

- 7.10. (Mixture of distributions.) Suppose that  $U$  is a  $Bernoulli(p(\theta))$  rv and conditionally, given  $U$ ,  $X$  is distributed with pdf  $f_U(x, \theta)$ . Then the joint distribution of  $(U, X)$  is

$$g(u, x, \theta) = \{p(\theta)\}^u \{1 - p(\theta)\}^{1-u} f_u(x, \theta), \quad u = 0 \text{ or } 1, x \in R, \theta \in \Theta \subset R,$$

where  $f_0(x, \theta)$  and  $f_1(x, \theta)$  are pdf's/pmf's on  $R$ . Assume that the usual regularity conditions hold for  $f_0, f_1$  and assume differentiability conditions on  $p(\theta)$  as needed. Let  $(U_1, X_1), \dots, (U_n, X_n)$  be iid, as  $(U, X)$  and let  $\hat{\theta}_n$  denote the MLE of  $\theta$

based on  $(U_i, X_i), i = 1, \dots, n$ . Allow for the possibility for  $f_0$  being discrete and  $f_1$  continuous.

**(a)** Show that

$$\begin{aligned} \text{E}_\theta \left[ \frac{\partial \log g(U, X; \theta)}{\partial \theta} \right] &= 0, \text{ and} \\ I(\theta) &= \text{E}_\theta \left[ -\frac{\partial^2 \log g(U, X, \theta)}{\partial \theta^2} \right] = \text{E}_\theta \left[ \left( \frac{\partial \log g(U, X, \theta)}{\partial \theta} \right)^2 \right] \\ &= \text{Var}_\theta \left( \frac{\partial \log g(U, X, \theta)}{\partial \theta} \right) \\ &= \frac{\{p'(\theta)\}^2}{p(\theta)(1-p(\theta))} + p(\theta)I_0(\theta) + (1-p(\theta))I_1(\theta), \quad \text{where} \\ I_u(\theta) &= \text{E}_\theta \left( \frac{\partial \log f_u(X, \theta)}{\partial \theta} \right)^2, \quad u = 0, 1. \end{aligned}$$

- (b)** How would you calculate the MLE  $\hat{\theta}_n$  of  $\theta$  and what is the asymptotic distribution of  $\hat{\theta}_n$ ?
- (c)** Let  $T_1, \dots, T_n$  be iid exponential rv's which are right-censored at  $c$ , giving rise to observations  $X_i = T_i I_{(0,c)}(T_i) + c I_{[c,\infty)}(T_i)$ . Here  $U_i = I_{(0,c)}(T_i)$  and conditionally on  $U_i$ ,  $X_i$  is degenerate at  $\{c\}$  if  $U_i = 0$  and had pdf  $(1/\theta)e^{-x/\theta}(1 - e^{-c/\theta})$  on  $(0, c)$  if  $U_i = 1$ . Discuss the maximum likelihood estimation of  $\theta$  from such data.
- 7.11.** Let  $X_1, \dots, X_n$  be iid with pdf/pmf  $f(x, \theta)$  and let  $\hat{\theta}_n$  denote the MLE of  $\theta$  based on  $X_i, i = 1, \dots, n$ . Under regularity conditions in [Section 7.1.1](#) show that the remainder terms in the expansion

$$0 = n^{-1/2} \sum_{i=1}^n \dot{l}(X_i, \theta) + \sqrt{n}(\hat{\theta}_n - \theta) \left[ n^{-1} \sum_{i=1}^n \ddot{l}(X_i, \theta) + n^{-1} \sum_{i=1}^n R_n(X_i) \right]$$

satisfies  $n^{-1} \sum_{i=1}^n R_n(X_i) = o_P(1)$ . Extend this result to the case when  $\theta$  is  $k$ -dim.

- 7.12.** Let  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  be independent samples from  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$ , respectively. Let  $m, n \rightarrow \infty$  so that  $m/(m+n) \rightarrow \alpha \in (0, 1)$ .
- (a)** Find the MLE of  $(\mu_1, \mu_2, \sigma^2)$  and its asymptotic distribution.
- (b)** Find the MLE of  $(\mu_1, \mu_2, \sigma^2)$  under the restriction  $\mu_1 = \mu_2$ .
- (c)** Derive the LRT statistic  $\Lambda_{m,n}$  for  $H_0: \mu_1 = \mu_2$  vs  $H_1: \mu_1 \neq \mu_2$ , reduce it to a suitable form and find the asymptotic distribution of  $-2 \log \Lambda_{m,n}$  under  $H_0$ ,
- (i)** using the general properties of LRT statistics, and
  - (ii)** from elementary considerations.

- 7.13.** Suppose that in [Section 7.1.1](#), the regularity condition

$$\begin{aligned} \sup_{|\theta - \theta_0| \leq \varepsilon} |\ddot{l}(x, \theta) - \ddot{l}(x, \theta_0)| &\leq H(x, \theta_0) \phi(\varepsilon) \text{ with} \\ \lim_{\varepsilon \rightarrow 0} \phi(\varepsilon) &= 0 \text{ and } \text{E}[H(X, \theta_0)] < \infty \end{aligned}$$

is strengthened by requiring

$$|\ddot{l}(x, \theta)| \leq H(x) \text{ with } E[H(X)] < \infty.$$

Show that the BAN property of

$$\tilde{\theta}_{n1} = \tilde{\theta}_{n0} - \sum \dot{l}(X_i, \tilde{\theta}_{n0}) / \sum \ddot{l}(X_i, \tilde{\theta}_{n0})$$

then holds if  $n^{1/4}(\tilde{\theta}_{n0} - \theta_0) = o_P(1)$ .

- 7.14.** Let  $(X_1, \dots, X_n)$  be a random sample from a distribution with pdf/pmf  $f(x, \theta)$  where  $\theta = (\theta_1, \dots, \theta_k)^T \in \Theta$  which is an open interval in  $\mathbb{R}^k$ , and  $\theta_0 = (\theta_{01}, \dots, \theta_{k0})^T$  is an interior point of  $\Theta$ . Suppose that the family  $\{f(x, \theta) : \theta \in \Theta\}$  satisfies the usual regularity conditions. We want to test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \in \Theta - \{\theta_0\}$ . Show that the following test statistics are equivalent (ie, differ from one another by  $o_P(1)$ ) under  $H_0$ :

- (a) Likelihood ratio statistic:  $T_{1n} = -2 \log \Lambda_n$ , where

$$\Lambda_n = \frac{\prod_{i=1}^n f(X_i, \theta_0)}{\left\{ \sup_{\theta \in \Theta} \prod_{i=1}^n f(X_i, \theta) \right\}},$$

- (b) Wald's statistic:  $T_{2n} = \mathbf{D}_n^T \mathbf{I}(\hat{\theta}_n) \mathbf{D}_n$ , where  $\mathbf{I}(\hat{\theta}_n)$  is the Fisher-information matrix evaluated at the MLE  $\hat{\theta}_n$  of  $\theta$  and  $\mathbf{D}_n = \sqrt{n}(\hat{\theta}_n - \theta_0)$ ,  
(c) Rao's statistic:  $T_{3n} = \mathbf{V}_n^T \mathbf{I}(\theta_0)^{-1} \mathbf{V}_n$ , where  $\mathbf{V}_n = n^{-1/2} \sum \dot{l}(X_i, \theta_0)$ .

- 7.15.** Let  $(X_{i1}, \dots, X_{in_i})$  be independent outcomes of Bernoulli trials in which  $X_{ij}$  takes values 1 or 0 with probabilities  $\theta_i$  and  $1 - \theta_i$ , respectively,  $i = 1, \dots, 4$ . We want to test  $H_0: \theta_{i+2}/\theta_{i+1} = \theta_{i+1}/\theta_i$ ,  $i = 1, 2$ , against the alternative  $H_1: \theta_{i+2}/\theta_{i+1} \neq \theta_{i+1}/\theta_i$  for at least one  $i$ . In this situation, computation of the LRT statistic  $-2 \log \Lambda_n$  becomes messy. The following approach based on Wald's statistic involves only the unrestricted MLEs  $\hat{\theta}_i = S_i/n_i$  of  $\theta_i$  where  $S_i = \sum_{1 \leq j \leq n_i} X_{ij}$ , but leads to a test criterion which is asymptotically equivalent to  $-2 \log \Lambda_n$  under  $H_0$  (see Exercise 7.14 above). For the asymptotics, let  $n_1, \dots, n_4 \rightarrow \infty$  in such a way that  $n_i/\sum n_k \rightarrow c_i > 0$ ,  $i = 1, \dots, 4$ . Since  $\theta_2^2 - \theta_1\theta_3 = \theta_3^2 - \theta_2\theta_4 = 0$  under  $H_0$ , large values of  $|\hat{\theta}_2^2 - \hat{\theta}_1\hat{\theta}_3|$  and  $|\hat{\theta}_3^2 - \hat{\theta}_2\hat{\theta}_4|$  would indicate departure from  $H_0$ .

- (a) Show that under  $H_0$ ,  $h(\hat{\theta})^T = \sqrt{n}(\hat{\theta}_2^2 - \hat{\theta}_1\hat{\theta}_3, \hat{\theta}_3^2 - \hat{\theta}_2\hat{\theta}_4)$  is asymptotically bivariate normal with mean vector  $\mathbf{0}$  and find the covariance matrix  $\Sigma(\theta)$  of this limiting distribution.  
(b) Find the asymptotic distribution of  $T_n = h(\hat{\theta})^T \Sigma(\hat{\theta})^{-1} h(\hat{\theta})$  under  $H_0$  and justify your answer.  
(c) Explain how you would find the critical value of a test for  $H_0$  based on  $T_n$  at level  $\alpha$ .

- 7.16.** Let  $(X_{i1}, \dots, X_{in_i})$  denote independent samples from  $Poi(\theta_i)$ ,  $i = 1, \dots, k$ . We want to test  $H_0: \theta_1 = \dots = \theta_k$  against all possible alternatives in  $(0, \infty)^k$ . Let  $T_i = \sum_{1 \leq j \leq n_i} X_{ij}$ ,  $\bar{T}_i = T_i/n_i$ ,  $n = \sum n_i$ ,  $T = \sum T_i$ , and  $\bar{T} = T/n$ , then  $(\bar{T}_1, \dots, \bar{T}_k)$  and  $(\bar{T}, \dots, \bar{T})$  are, respectively, the unrestricted MLE and the restricted MLE of

$(\theta_1, \dots, \theta_k)$  under  $H_0$ . [ $T = 0$  causes some difficulties which we shall ignore for large sample purposes.]

- (a) Show that the conditional distribution of  $(T_1, \dots, T_k)$  given  $T$  is  $Multi(T, n_1/n, \dots, n_k/n)$ .
- (b) Justify the use of  $W_n = \sum[(T_i - n_i \bar{T})^2 / (n_i \bar{T})]$  as a test statistic for testing  $H_0$ . What should be the rejection region of such a test at level  $\alpha$  when  $n$  is large?
- (c) Let  $\Lambda_n$  denote the LRT statistic for  $H_0$ . Show that

$$-2 \log \Lambda_n = 2 \sum n_i \bar{T}_i \log(\bar{T}_i / \bar{T}).$$

What is the asymptotic distribution of  $-2 \log \Lambda_n$  under  $H_0$ ?

- (d) Show that under  $H_0$ ,  $-2 \log \Lambda_n = W_n + o_P(1)$ . [Notice the similarity with Example 7.2.1 at the end of Section 7.2.]

- 7.17. The number of flies ( $X_1, X_2, X_3$ ) in three categories, resulting from certain crossings is distributed as  $Multi(n, \pi)$ . According to the Hardy-Weinberg formula, the probabilities of this multinomial distribution are

$$\pi_1(\theta) = (1 - \theta)^2, \pi_2(\theta) = 2\theta(1 - \theta), \pi_3(\theta) = \theta^2,$$

for some  $0 < \theta < 1$ . In an experiment, the observed frequencies in the three categories are  $x_1 = 45$ ,  $x_2 = 58$ , and  $x_3 = 22$ , in a random sample of size  $n = 125$ . Test whether the data support the above model at a level of significance  $\alpha = 0.05$ .

# Distribution-Free Tests for Hypothesis Testing in Nonparametric Families

## 8.1 Ranks and Order Statistics

We start with the general case in which  $\mathbf{X} = (X_1, \dots, X_n)^T$  is a random vector with the joint pdf  $f(\mathbf{x})$  on  $\mathbb{R}^n$ , and in this setting derive the joint distribution of the vector of order statistics  $\mathbf{X}_{(n)} = (X_{n:1}, \dots, X_{n:n})^T$  where  $X_{n:1} < \dots < X_{n:n}$  and the vector of ranks  $\mathbf{R}_n = (R_{n:1}, \dots, R_{n:n})^T$  where  $R_{n:i} = 1 + \sum_{j \neq i}^n I_{(0,\infty)}(X_i - X_j)$  is the rank of  $X_i$  among  $X_1, \dots, X_n$ .

**Theorem 8.1.1.** *If  $(X_1, \dots, X_n)$  has joint pdff on  $\mathbb{R}^n$ , then*

(i) *the pdf of  $\mathbf{X}_{(n)}$  is*

$$\bar{f}(y_1, \dots, y_n) = \sum_{\mathbf{r}} f(y_{r_1}, \dots, y_{r_n}), \quad y_1 < \dots < y_n,$$

*where the sum is over all  $n!$  permutations  $\mathbf{r} = (r_1, \dots, r_n)^T$  of  $(1, \dots, n)$ , and*

(ii)  $P[\mathbf{R} = \mathbf{r} | \mathbf{X}_{(n)} = \mathbf{y}] = f(y_{r_1}, \dots, y_{r_n}) / \bar{f}(y_1, \dots, y_n)$ .

*Proof.* For  $A \subset \{\mathbf{y}: y_1 < \dots < y_n\}$ ,

$$P[\mathbf{X}_{(n)} \in A] = \int_{\mathbf{x}_{(n)} \in A} f(\mathbf{x}) d\mathbf{x} = \sum_{\mathbf{r}} \int_{r(\mathbf{x}) = \mathbf{r}, \mathbf{x}_{(n)} \in A} f(\mathbf{x}) d\mathbf{x},$$

where it is understood that  $\int$  in the above expressions refer to an  $n$ -dimensional integral and  $d\mathbf{x} = dx_1 \cdots dx_n$ . On  $\{\mathbf{x}: r(\mathbf{x}) = \mathbf{r}\}$ , transform  $(y_1, \dots, y_n) = (x_{n:1}, \dots, x_{n:n})$  which is one-to-one with Jacobian equal to 1, and  $r(\mathbf{x}) = \mathbf{r} \iff y_{r_i} = x_{n:r_i} = x_i$ . Thus

$$P[\mathbf{X}_{(n)} \in A] = \sum_{\mathbf{r}} \int_A f(y_{r_1}, \dots, y_{r_n}) d\mathbf{y} = \int_A \left\{ \sum_{\mathbf{r}} f(y_{r_1}, \dots, y_{r_n}) \right\} d\mathbf{y},$$

which proves (i). Moreover,

$$\begin{aligned} P[\mathbf{R} = \mathbf{r}, \mathbf{X}_{(n)} \in A] &= \int_{r(\mathbf{x}) = \mathbf{r}, \mathbf{x}_{(n)} \in A} f(\mathbf{x}) d\mathbf{x} = \int_A f(y_{r_1}, \dots, y_{r_n}) d\mathbf{y} \\ &= \int_A \frac{f(y_{r_1}, \dots, y_{r_n})}{\bar{f}(y_1, \dots, y_n)} \bar{f}(y_1, \dots, y_n) d\mathbf{y}, \end{aligned}$$

proving (ii). □

**Corollary 8.1.1.** If  $f(y_{r_1}, \dots, y_{r_n}) = f(y_1, \dots, y_r)$  for all  $r$ , then

- (i) the pdf of  $\mathbf{X}_{(n)}$  equals  $n!f(y_1, \dots, y_n)$ ,  $y_1 < \dots < y_n$ ,
- (ii)  $P[\mathbf{R} = \mathbf{r} | \mathbf{X}_{(n)} = \mathbf{y}] = 1/n!$  for all  $\mathbf{r}$  and  $\mathbf{y}$ , so  $P[\mathbf{R} = \mathbf{r}] = 1/n!$  for all  $\mathbf{r}$ , and  $\mathbf{X}_{(n)}$  and  $\mathbf{R}$  are mutually independent.

In particular, if  $X_1, \dots, X_n$  are iid with common pdf  $f$ , then the pdf of  $\mathbf{X}_{(n)}$  is  $n! \prod_{i=1}^n f(y_i)$ ,  $y_1 < \dots < y_n$  as already proved in [Section 1.11](#).

For the next result, restrict to  $X_1, \dots, X_n$  which are iid with common pdf  $f$ . With slight abuse of notation, we shall write the joint pdf of  $(X_1, \dots, X_n)$  as  $f(\mathbf{x}) = \prod_{i=1}^n f(x_i)$ .

Suppose that the common pdf of  $f$  of  $X_1, \dots, X_n$  is symmetric (ie,  $f(-x) = f(x)$  for all  $x$ ). Consider

$$\mathbf{s} = (\text{sign}(X_1), \dots, \text{sign}(X_n))^T, |\mathbf{X}| = (|X_1|, \dots, |X_n|)^T, \text{ and let} \\ |\mathbf{X}|_{(n)} = (|X|_{n:1}, \dots, |X|_{n:n})^T$$

denote the vector of order statistics and  $\mathbf{R}_n^+ = (R_{n:1}^+, \dots, R_{n:n}^+)$  denote the ranks of  $|X_i|$  among  $|X_1|, \dots, |X_n|$ . Then

- (i) for each  $i$ ,  $|X_i|$  and  $\text{sign}(X_i)$  are independent,

$$P[S_i = 1] = P[S_i = -1] = 1/2,$$

and  $|X_i|$  has pdf  $2f(x)$ ,  $x > 0$ .

- (ii) since  $X_1, \dots, X_n$  are independent,  $S_1, \dots, S_n$  and  $|X_1|, \dots, |X_n|$  are all mutually independent,
- (iii) also, the ranks  $\mathbf{R}_n^+$  and the order statistics  $|\mathbf{X}|_{(n)}$  are independent, having distributions obtained in the Corollary of [Theorem 8.1.1](#).

We thus have

**Theorem 8.1.2.** If  $X_1, \dots, X_n$  are iid with common pdf  $f$  which is symmetric about 0, then

- (i) the vectors  $\mathbf{S}$ ,  $\mathbf{R}_n^+$ , and  $|\mathbf{X}|_{(n)}$  are mutually independent,
- (ii)  $P[\mathbf{R}_n^+ = \mathbf{r}] = 1/n!$  for all  $\mathbf{r}$ ,
- (iii)  $|\mathbf{X}|_{(n)}$  has joint pdf  $n!2^n \prod_{i=1}^n f(y_i)$ ,  $0 < y_1 < \dots < y_n$ ,
- (iv)  $P[\mathbf{S} = \mathbf{s}] = 1/2^n$  for all  $\mathbf{s} = (\pm 1, \dots, \pm 1)$ .

The following rank-related lemma is for future use.

**Lemma 8.1.1.** Let  $t(X_1, \dots, X_n)$  be a function of iid rv's  $X_1, \dots, X_n$ , and suppose  $t(\mathbf{X})$  has finite expectation. Then

$$E[t(X_1, \dots, X_n) | \mathbf{R} = \mathbf{r}] = E[t(X_{n:r_1}, \dots, X_{n:r_n}) | \mathbf{R}_n = \mathbf{r}] \\ = E[t(X_{n:r_1}, \dots, X_{n:r_n})],$$

because  $\mathbf{X}_{(n)}$  and  $\mathbf{R}_n$  are independent.

## Application: Permutation Test

Let  $\mathcal{P}_0$  be a family of pdf's on  $\mathbb{R}^n$  defined by

$$\mathcal{P}_0 = \left\{ p: p(x_1, \dots, x_N) = \prod_{i=1}^N f(x_i) \quad \text{where } f \text{ is an arbitrary pdf on } \mathbb{R} \right\}.$$

On the basis of a random sample from a distribution with pdf  $p$ , we want to test  $H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q \notin \mathcal{P}_0$  for a specified  $q$ .

Since  $(X_1, \dots, X_N)$  is in one-to-one correspondence with  $(\mathbf{X}_{(N)}, \mathbf{R}_N)$ , the vectors of order statistics and ranks, we consider tests  $\Psi(\mathbf{X}_{(N)}, \mathbf{R}_{(N)})$ .

From Chapter 6, recall that a test  $\Psi$  is a similar region test of size  $\alpha$  for  $H_0: p \in \mathcal{P}_0$  if

$$E_p[\Psi(\mathbf{X}_N, \mathbf{R}_N)] = \alpha \text{ for all } p \in \mathcal{P}_0.$$

Since the conditional distribution of  $(X_1, \dots, X_N)$  given  $\mathbf{X}_{(N)}$  is uniformly distributed over the set of  $N!$  permutations of  $(X_1, \dots, X_N)$  irrespective of  $p \in \mathcal{P}_0$ ,  $\mathbf{X}_{(N)}$  is sufficient for  $p \in \mathcal{P}_0$  in  $(X_1, \dots, X_N)$ . It can also be shown that  $\mathbf{X}_{(N)}$  is complete, that is,  $E_p[g(\mathbf{X}_{(N)})] = 0$  for all  $p \in \mathcal{P}_0$  implies  $g(\mathbf{x}_{(N)}) = 0$ , a.s.  $\mathcal{P}_0$ . By virtue of complete sufficiency of  $\mathbf{X}_{(N)}$ , the similar region property of  $\Psi$  holds iff  $\Psi$  has Neyman-structure with respect to  $\mathbf{X}_{(N)}$  by Theorem 6.9.1, ie, for almost all  $\mathbf{x}_{(N)}$ ,

$$E_{H_0}[\Psi(\mathbf{X}_{(N)}, \mathbf{R}_N) | \mathbf{X}_{(N)} = \mathbf{x}_{(N)}] = \alpha.$$

Moreover, since  $\mathbf{X}_{(N)}$  and  $\mathbf{R}_N$  are independent under  $H_0$ , the last expression can be rewritten as

$$(1/N!) \sum_r \Psi(\mathbf{x}_{(N)}, \mathbf{r}) = \alpha \text{ for all almost } \mathbf{x}_{(N)}, \quad (1)$$

and the problem of finding the most powerful test at level  $\alpha$  for  $H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q \notin \mathcal{P}_0$ ; that is, the problem of maximizing  $E_q[\Psi(\mathbf{X}_{(N)}, \mathbf{R}_N)]$  subject to Eq. (1) is solved by maximizing

$$E_q[\Psi(\mathbf{X}_{(N)}, \mathbf{R}_N) | \mathbf{X}_{(N)} = \mathbf{x}_{(N)}] = E_q[\Psi(\mathbf{x}_{(N)}, \mathbf{R}_N) | \mathbf{X}_{(N)} = \mathbf{x}_{(N)}]$$

subject to Eq. (1) for each  $\mathbf{x}_{(N)}$ .

The optimal  $\Psi$  is obtained by using the N-P Lemma conditionally, given  $\mathbf{X}_{(N)} = \mathbf{x}_{(N)}$ . By Theorem 8.1.1 and its Corollary, the conditional likelihood ratio

$$P_q[\mathbf{R}_N = \mathbf{r}_N | \mathbf{X}_{(N)} = \mathbf{x}_{(N)}] / P_p[\mathbf{R}_N = \mathbf{r}_N | \mathbf{X}_{(N)} = \mathbf{x}_{(N)}]$$

can be equivalently expressed as

$$N! q(x_{N:r_1}, \dots, x_{N:r_N}) / q(\mathbf{x}_{(N)}).$$

The optimal  $\Psi$  is therefore given by

$$\begin{aligned} \Psi(\mathbf{x}_{(N)}, \mathbf{r}_N) &= 0, \text{ or } \gamma(\mathbf{x}_{(N)}), \text{ or } 1, \text{ according as} \\ q(x_{N:r_1}, \dots, x_{N:r_N}) <, \text{ or } =, \text{ or } > k(\mathbf{x}_{(N)}), \end{aligned}$$

where  $k(\mathbf{x}_{(N)})$  and  $0 \leq \gamma(\mathbf{x}_{(N)}) \leq 1$  are determined by the size  $\alpha$  condition.

To find  $k(\mathbf{x}_{(N)})$  and  $\gamma(\mathbf{x}_{(N)})$  for a given  $\mathbf{x}_{(N)}$ , we arrange the  $N!$  values of

$$q_0(\mathbf{x}_{(N)}, \mathbf{r}_N) = q(x_{N:r_1}, \dots, x_{N:r_N})$$

for all permutations  $\mathbf{r}_N$  of  $(1, \dots, N)$  in ascending order of magnitude and let  $\mathbf{r}_N^*$  be such that

$$\begin{aligned} v_N &= \#\{\mathbf{r}_N : q_0(\mathbf{x}_{(N)}, \mathbf{r}_N) \geq q_0(\mathbf{x}_{(N)}, \mathbf{r}_N^*)\} \geq N!\alpha, \text{ but} \\ v_N^- &= \#\{\mathbf{r}_N : q_0(\mathbf{x}_{(N)}, \mathbf{r}_N) > q_0(\mathbf{x}_{(N)}, \mathbf{r}_N^*)\} < N!\alpha. \end{aligned}$$

Then

$$k(\mathbf{x}_{(N)}) = q_0(\mathbf{x}_{(N)}, \mathbf{r}_N^*) \text{ and } \gamma(\mathbf{x}_{(N)}) = (N!\alpha - v_N^-)/(v_N - v_N^-).$$

**Example 8.1.1.** For  $N = 4$ , let  $q(x_1, x_2, x_3, x_4) = f(x_1)f(x_2)f(x_3)g(x_4)$ , where  $f(x) = 1$ ,  $0 \leq x \leq 1$  and  $g(x) = 2x$ ,  $0 \leq x \leq 1$ , and let  $p$  denote the unknown joint pdf of  $(X_1, \dots, X_4)$ . We want to test  $H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q$  at level  $\alpha = 0.1$  based on observations  $(x_1, \dots, x_4) = (0.2, 0.7, 0.4, 0.9)$ .

*Solution.* The vectors of ranks and order statistics in the observed data are  $\mathbf{r}_4 = (1, 3, 2, 4)$  and  $\mathbf{x}_{(4)} = (0.2, 0.4, 0.7, 0.9)$ . The values of  $q(x_{4:r_1}, \dots, x_{4:r_4})$  are

$$\begin{aligned} 4x_3x_4 &= 4(0.4)(0.9) = 1.44 \text{ for } \mathbf{r}_4 = (1, 2, 3, 4), (1, 2, 4, 3), (2, 1, 3, 4), (2, 1, 4, 3), \\ 4x_1x_2 &= 4(0.2)(0.7) = 0.56 \text{ for } \mathbf{r}_4 = (3, 4, 1, 2), (3, 4, 2, 1), (4, 3, 1, 2), (4, 3, 2, 1), \\ 4x_2x_4 &= 4(0.7)(0.9) = 2.52 \text{ for } \mathbf{r}_4 = (1, 3, 2, 4), (1, 3, 4, 2), (3, 1, 2, 4), (3, 1, 4, 2), \\ 4x_1x_3 &= 4(0.2)(0.4) = 0.32 \text{ for } \mathbf{r}_4 = (2, 4, 1, 3), (2, 4, 3, 1), (4, 2, 1, 3), (4, 2, 3, 1), \\ 4x_2x_3 &= 4(0.7)(0.4) = 1.12 \text{ for } \mathbf{r}_4 = (1, 4, 2, 3), (1, 4, 3, 2), (4, 1, 2, 3), (4, 1, 3, 2), \\ 4x_1x_4 &= 4(0.2)(0.9) = 0.72 \text{ for } \mathbf{r}_4 = (2, 3, 1, 4), (2, 3, 4, 1), (3, 2, 1, 4), (3, 2, 4, 1). \end{aligned}$$

These values of  $q(\cdot)$  are arranged as

$$0.32 < 0.56 < 0.72 < 1.12 < 1.44 < 2.52,$$

each value repeated four times, and  $N!\alpha = (24)(0.1) = 2.40$ .

Hence

$$\begin{aligned} q_0(\mathbf{x}_{(4)}, \mathbf{r}_4^*) &= 2.52, v_4 = 4, \text{ and } v_4^- = 0, \text{ so} \\ k_n &= 2.52 \text{ and } \gamma = 2.40/4 = 0.60. \end{aligned}$$

For the observed data,  $\mathbf{r}_4 = (1, 3, 2, 4)$  and  $q_0 = 2.52$ , so we randomize and reject  $H_0$  with probability  $\gamma = 0.60$ .

*Remark 8.1.1.* The main problem in implementing a permutation test is that we cannot use standard tables. The critical value of  $q_0(\mathbf{x}_{(N)}, \mathbf{r}_N)$  must be determined in every instance by the observed data.

### 8.1.1 Nonparametric Tests in Three Basic Problems

#### 1. Test for symmetry in the one-sample problem.

Let  $X_1, \dots, X_n$  be iid with common continuous cdf  $F$  (Unknown). We want to test  $H_0: F(-x) = 1 - F(x)$  for all  $x$  (ie, the distribution is symmetric about zero), vs  $H_1$  the

distribution is not symmetric about zero, or vs  $H_+$ :  $F(-x) \leq 1 - F(x)$  for all  $x$  with strict inequality for some  $x$  (or  $\int_{-\infty}^{\infty} [1 - F(-x)] dF(x) > 1/2$ ).

Let  $S_i = \text{sign}(X_i)$  and  $R_{n:i}^+ = \text{rank of } |X_i| \text{ among } |X_1|, \dots, |X_n|$ .

Under  $H_0$ , high ranks and low ranks among  $\{R_{n:i}^+\}$  will be equally associated with  $S_i = 1$  or  $-1$ , but under  $H_+$ , the observations with  $S_i = 1$  will tend to have higher  $R_{n:i}^+$  and those with  $S_i = -1$  will tend to have lower  $R_{n:i}^+$ . This leads to the consideration of the Wilcoxon signed-rank statistic:

$$T_n = \{n(n-1)\}^{-1} \sum_{i=1}^n S_i R_{n:i}^+ \text{ or equivalently, } \sum_{\{i: X_i > 0\}} R_{n:i}^+ = (n + T_n)/2.$$

We reject  $H_0$  in favor of  $H_+$  (or  $H_1$ ) if  $T_n$  (or  $|T_n|$ )  $\geq$  critical value.

## 2. Test for homogeneity in the two-sample problem.

Let  $X_1, \dots, X_m$  be iid with common continuous cdf  $F$  and let  $Y_1, \dots, Y_n$  be iid with common continuous cdf  $G$ , the two samples being mutually independent. We want to test

$$\begin{aligned} H_0 : F = G \text{ vs } H_1: P_{F,G}[X_i > Y_j] &= \int_{-\infty}^{\infty} G(x) dF(x) \neq 1/2, \text{ or} \\ H_+: P_{F,G}[X_i > Y_j] &> 1/2. \end{aligned}$$

For notational convenience, write the combined sample as

$$X_1, \dots, X_m, Y_1 = X_{m+1}, \dots, Y_n = X_{m+n}.$$

The average ranks of the  $X$ -observations and the  $Y$ -observations in the combined sample are  $m^{-1} \sum_{i=1}^m R_{m+n:i}$  and  $n^{-1} \sum_{i=m+1}^{m+n} R_{m+n:i}$ , respectively. Since the average ranks of the two samples would tend to be equal under  $H_0$ , the difference between the average ranks, known as the Wilcoxon two-sample rank-sum statistic defined as

$$\begin{aligned} W_{m,n} &= m^{-1} \sum_{i=1}^m R_{m+n:i} - n^{-1} \sum_{i=m+1}^{m+n} R_{m+n:i}, \text{ or equivalently} \\ W'_{m,n} &= \sum_{i=1}^m R_{m+n:i} \text{ (which is a linear function of } W_{m,n}) \end{aligned}$$

is an indicator of departure from  $H_0$ . Moreover, by algebraic rearrangement, we can write

$$W_{m,n} = (m+n)(T_{m,n} - 1/2), \text{ where } T_{m,n} = (mn)^{-1} \sum_{i=1}^m \sum_{j=1}^n I_{(0,\infty)}(X_i - Y_j)$$

which is called the Mann-Whitney statistic. Thus the three statistics  $W_{m,n}$ ,  $W'_{m,n}$ , and  $T_{m,n}$  are linear functions of one another and any of them can be used as a test statistic for testing  $H_0$  vs  $H_1$  or  $H_0$  vs  $H_+$ . In particular, with the Mann-Whitney statistic, we can reject  $H_0$  if favor of  $H_+$  (or  $H_1$ ) if  $T_{m,n}$  (or  $|T_{m,n}|$ )  $\geq$  critical value.

### 3. Test for independence in a bivariate distribution.

Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be iid as  $(X, Y)$ . We want to test  $H_0$ :  $X$  and  $Y$  are independent vs  $H_1$ :  $X$  and  $Y$  are dependent, or  $H_+$ :  $X$  and  $Y$  are positively dependent. Let  $U, V, Z$  be independent rv's with unknown pdf's  $f, g, h$ . Then the alternative hypothesis can be formulated as  $(X, Y) = (U + \Delta Z, V + \Delta Z)$ , and then take  $H_1$ :  $\Delta \neq 0$  and  $H_+$ :  $\Delta > 0$ .

Let  $R_{n:i}$  = rank of  $X_i$  among  $X_1, \dots, X_n$  and  $R'_{n:i}$  = rank of  $Y_i$  among  $Y_1, \dots, Y_n$ . Then replacing  $(X_1, Y_1), \dots, (X_n, Y_n)$  by their ranks  $(R_{n:1}, R'_{n:1}), \dots, (R_{n:n}, R'_{n:n})$ , the Spearman's rank correlation defined as

$$\begin{aligned}\rho_S &= \text{correlation coefficient between } \{(R_{n:i}, R'_{n:i}), i = 1, \dots, n\} \\ &= 12 \left\{ n(n^2 - 1) \right\}^{-1} \sum_{i=1}^n R_{n:i} R'_{n:i} - 3(n+1)/(n-1)\end{aligned}$$

(using  $\bar{R}_n = \bar{R}'_n = (n+1)/2$  and  $\sum_{i=1}^n (R_{n:i} - \bar{R}_n)^2 = \sum_{i=1}^n (R'_{n:i} - \bar{R}'_n)^2 = n(n^2 - 1)/12$ ) is an obvious candidate for testing  $H_0$  vs  $H_1$  or  $H_+$ . One can use either  $\rho_S$  or  $\sum_{i=1}^n R_{n:i} R'_{n:i}$  as a test statistic.

Another test statistic for this problem is based on the simple idea that if we compare  $(X_i, Y_i)$  with  $(X_j, Y_j)$ , then  $X_i - X_j$  being positive or negative is independent of  $Y_i - Y_j$  being positive or negative if  $X, Y$  are independent. An overall comparison between all pairs of data-points leads to

$$\tau_n = [n(n-1)]^{-1} \sum_{i=1}^n \sum_{j \neq i, j=1}^n \text{sign}(X_i - X_j) \text{sign}(Y_i - Y_j),$$

known as Kendall's tau statistic. We reject  $H_0$  in favor of  $H_+$  (or  $H_1$ ) if  $\tau_n$  (or  $|\tau_n|$ )  $\geq$  critical value.

The statistic  $\rho_S$  differs from a linear function of  $\tau_n$  by  $o_P(1)$  as  $n \rightarrow \infty$  under  $H_0$ , so a test based on  $\rho_S$  is carried out analogously (see [4]).

#### 8.1.2 Exact Distribution of Rank Statistic Under $H_0$

Exact distributions of rank statistics under  $H_0$  can be obtained by combinatorial arguments, using the properties of ranks under  $H_0$ . We illustrate this with the example of Mann-Whitney statistic in the two-sample problem.

Under  $H_0$ , the rv's  $X_1, \dots, X_m, Y_1, \dots, Y_n$  are iid. Let  $v_{m,n}(u)$  be the number of arrangements of  $m$   $X$ 's and  $n$   $Y$ 's in which exactly  $u$  pairs have  $X > Y$  (ie,  $I_{(0,\infty)}(X - Y) = u$ ). Since these arrangements are equally likely under  $H_0$ , we have

$$P_{H_0} \left[ \sum_{i=1}^m \sum_{j=1}^n I_{(0,\infty)}(X_i - Y_j) = u \right] = v_{m,n}(u) / \binom{m+n}{m}.$$

The main thing is to find a formula for  $v_{m,n}(u)$ , which is obtained from the difference equation

$$v_{m,n}(u) = v_{m,n-1}(u) + v_{m-1,n}(u-n),$$

where  $v_{m,n-1}(u)$  corresponds to sequences ending with  $Y$  and  $v_{m-1,n}(u-n)$  corresponds to sequences ending in  $X$ . Using this difference equation with initial conditions

$$v_{m,0}(u) = \begin{cases} 1 & u=0 \\ 0 & u \neq 0 \end{cases} \text{ and } v_{0,n}(u) = \begin{cases} 1 & u=0 \\ 0 & u \neq 0 \end{cases},$$

we can compute  $v_{m,n}(u)$  recursively.

Tables of exact tail probabilities are available for this and many other rank statistics for small to moderate sample sizes.

### 8.1.3 Asymptotic Distribution of Rank Statistics Under $H_0$ by $U$ -Statistic Approach

*One-sample  $U$ -statistic.* Let  $X_1, \dots, X_n$  be iid with a common continuous cdf  $F$ . Many one-sample rank statistics are of the form

$$U_n = \left(1/n^{(r)}\right) \sum_{n,r} g(X_{i_1}, \dots, X_{i_r}),$$

where  $\sum_{n,r}$  denotes sum over all distinct  $i_1, \dots, i_r$  in  $\{1, \dots, n\}$  and  $n^{(r)} = n(n-1) \cdots (n-r+1)$ .

Assume without loss of generality that  $g$  is symmetric in its coordinates. [If not, replace  $g$  by  $g^*$  obtained by averaging over all permutations of its coordinates, which leaves  $U_n$  unchanged.] A statistic of this form is called  $U$ -statistic [33].

Suppose that  $E_F[g^2(X_1, \dots, X_r)] < \infty$ , and let

$$\begin{aligned} \theta &= \theta(F) = E_F[g(X_1, \dots, X_r)], \quad h(X_1, \dots, X_r) = g(X_1, \dots, X_r) - \theta \text{ and} \\ h_c(X_1, \dots, X_c) &= E_F[h(X_1, \dots, X_c, X_{c+1}, \dots, X_r) | X_1, \dots, X_c] \text{ for } c = 1, \dots, r, \text{ and} \\ \xi_c &= \xi_c(F) = E_F[h^2(X_1, \dots, X_c)]. \end{aligned}$$

Then the mean and variance of  $U_n$  are

$$E[U_n] = \theta \text{ and } \text{Var}[U_n] = \left(n^{(r)}\right)^{-2} \sum_{c=1}^r N_c \xi_c, \text{ where}$$

$$N_n = \#\{(i_1, \dots, i_r), (j_1, \dots, j_r) \text{ with exactly } c \text{ elements in common}\}$$

$$= \binom{n}{2r-c} \frac{(2r-c)!}{c!(r-c)!(r-c)!} (r!)^2 = \binom{n}{r} \binom{r}{c} \binom{n-r}{r-c} (r!)^2.$$

Suppose that  $\xi_1 > 0$ . Since  $N_{c+1}/N_c \rightarrow 0$  as  $n \rightarrow \infty$ ,  $\lim_{n \rightarrow \infty} n \text{Var}[U_n] = r^2 \xi_1 > 0$ .

Actually,  $U_n$  is asymptotically normal.

**Theorem 8.1.3.** If  $\xi_1 > 0$ , then  $\sqrt{n}(U_n - \theta) \xrightarrow{\mathcal{L}} Z \sim N(0, r^2 \xi_1)$ .

Postponing the proof of [Theorem 8.1.3](#), we look at *Two-sample U-statistics*. Let  $\mathbf{X} = (X_1, \dots, X_m)$  and  $\mathbf{Y} = (Y_1, \dots, Y_n)$  be mutually independent random samples from distributions with cdf's  $F$  and  $G$ , respectively, and define

$$U_{m,n} = \left\{ m^{(r)} n^{(s)} \right\}^{-1} \sum_{m,r} \sum_{n,s} g(X_{i_1}, \dots, X_{i_r}, Y_{j_1}, \dots, Y_{j_s}), \text{ where}$$

$$g(x_{i_1}, \dots, x_{i_r}, y_{j_1}, \dots, y_{j_s}) = g(x_1, \dots, x_r, y_1, \dots, y_s)$$

for all permutations  $(i_1, \dots, i_r)$  of  $(1, \dots, r)$  and  $(j_1, \dots, j_s)$  of  $(1, \dots, s)$ .

Suppose that  $E[g^2(X_1, \dots, X_r, Y_1, \dots, Y_s)] < \infty$  and define

$$\theta = \theta(F, G) = E_{F,G}[g(\mathbf{X}, \mathbf{Y})], \quad h(\mathbf{X}, \mathbf{Y}) = g(\mathbf{X}, \mathbf{Y}) - \theta$$

and for  $c = 0, 1, \dots, r$  and  $d = 0, 1, \dots, s$ , let

$$\xi_{c,d} = E[h_{c,d}^2(X_1, \dots, X_c, Y_1, \dots, Y_d)],$$

where

$$h_{c,d}(x_1, \dots, x_c, y_1, \dots, y_d) = E[h(\mathbf{X}, \mathbf{Y}) | X_1, \dots, X_c, Y_1, \dots, Y_d]$$

if  $\min(c, d) \geq 1$  and  $h_{0,0} = 0$ . The mean and variance of  $U_{m,n}$  are

$$\begin{aligned} E[U_{m,n}] &= \theta, \text{ and} \\ \text{Var}[U_{m,n}] &= \left\{ \binom{m}{r} \binom{n}{s} \right\}^{-2} \sum_{c=0}^r \sum_{d=0}^s \binom{r}{c} \binom{m-r}{r-c} \binom{s}{d} \binom{n-s}{s-d} \xi_{c,d} \\ &= (r^2/m) \xi_{1,0} + (s^2/n) \xi_{0,1} + o(1/\min(m, n)) \text{ as } m, n \rightarrow \infty. \end{aligned}$$

Thus if  $m, n \rightarrow \infty$  in such a way that  $m/(m+n) \rightarrow \lambda \in (0, 1)$ , then

$$\lim_{m,n \rightarrow \infty} (m+n) \text{Var}[U_{m,n}] = r^2 \xi_{1,0}/\lambda + s^2 \xi_{0,1}/(1-\lambda),$$

and  $U_{m,n}$  is asymptotically normal.

**Theorem 8.1.4.** *Under the above conditions,*

$$\sqrt{m+n}(U_{m,n} - \theta) \xrightarrow{\mathcal{L}} Z \sim N\left(0, r^2 \xi_{1,0}/\lambda + s^2 \xi_{0,1}/(1-\lambda)\right).$$

The proofs of [Theorems 8.1.3](#) and [8.1.4](#) are accomplished by the *Hájek Projection Method* of approximating an arbitrary function of independent rv's by a sum of functions of individual rv's.

Let  $Z_1, \dots, Z_n$  be independent rv's. Consider all rv's  $\varphi(Z_1, \dots, Z_n)$  with  $E[\varphi^2(Z)] < \infty$  as a vector space  $\mathcal{J}$  with inner product  $\langle T_1, T_2 \rangle = \text{Cov}[T_1, T_2]$ . In  $\mathcal{J}$ , let  $\mathcal{J}_0$  denote the subspace of all rv's of the form  $S = \sum_{i=1}^n \psi_i(Z_i)$ . We want to approximate an arbitrary  $T \in \mathcal{J}$  with  $E[T] = 0$  by an rv  $S \in \mathcal{J}_0$  with the smallest  $E[(T-S)^2] = \|T-S\|^2$  (ie, the projection of  $T$  on  $\mathcal{J}_0$ ).

**Theorem 8.1.5** (Hájek Projection Theorem). *If  $E[T] = 0$  and  $E[T^2] < \infty$ , then the projection of  $T$  on  $\mathcal{J}_0$  is  $T^* = \sum_{i=1}^n E[T|Z_i]$  w.p. 1.*

*Proof.* Let  $\psi_i^*(Z_i) = E[T|Z_i]$  and consider an arbitrary  $S = \sum_{i=1}^n \psi_i(Z_i) \in \mathcal{J}_0$ . To prove the theorem, it is enough to show that  $E[(T - T^*)(T^* - S)] = 0$ , because then we would have

$$E[(T - S)^2] = E[(T - T^*)^2] + E[(T^* - S)^2] \geq E[(T - T^*)^2], \quad (2)$$

with equality holding iff  $S = T^*$  w.p. 1. Now

$$\begin{aligned} E[(S - T^*)(T - T^*)] &= \sum_{i=1}^n EE \left[ \{\psi_i(Z_i) - \psi_i^*(Z_i)\} \left\{ T - \psi_i^*(Z_i) - \sum_{j \neq i, j=1}^n \psi_j^*(Z_j) \right\} \middle| Z_i \right] \\ &= \sum_{i=1}^n E \left[ \{\psi_i(Z_i) - \psi_i^*(Z_i)\} \left\{ E(T|Z_i) - \psi_i^*(Z_i) - \sum_{j \neq i, j=1}^n E(\psi_j^*(Z_j)) \right\} \right] \\ &= \sum_{i=1}^n E[\{\psi_i(Z_i) - \psi_i^*(Z_i)\}\{\psi_i^*(Z_i) - \psi_i^*(Z_i) - 0\}] = 0, \end{aligned}$$

since  $Z_j$  is independent of  $Z_i$  for  $j \neq i$ , and

$$E[\psi_j^*(Z_j)] = EE[T|Z_j] = E[T] = 0.$$

□

### Corollary 8.1.2.

(a) Taking  $S = 0$  in Eq. (2), we have

$$E[(T - T^*)^2] = E[T^2] - E[T^{*2}], \text{ ie, } Var[T - T^*] = Var[T] - Var[T^*],$$

(b) if  $E[T] = \mu$ , then  $T' = T - \mu$  has mean 0 and the projection of  $T'$  on  $\mathcal{J}_0$  is

$$\begin{aligned} \sum_{i=1}^n E[T'|Z_i] &= \sum_{i=1}^n E[T|Z_i] - n\mu. \text{ Hence the projection of } T = T' + \mu \text{ on } \mathcal{J}_0 \text{ is} \\ T^* &= \sum_{i=1}^n E[T|Z_i] - (n-1)E[T]. \end{aligned}$$

Proofs of Theorems 8.1.3 and 8.1.4.

We shall find the Hajek projections of the  $U$ -statistics and then apply the CLT.

In the one-sample problem, the Hajek projection of  $T_n = \sqrt{n}(U_n - \theta)$  is

$$\begin{aligned} T_n^* &= \sum_{i=1}^n E[T_n|X_i] = \sqrt{n} \{n^{(r)}\}^{-1} \sum_{i=1}^n E \left[ \sum_{n,r} h(X_{i_1}, \dots, X_{i_r}) | X_i \right] \\ &= rn^{-1/2} \sum_{i=1}^n h_1(X_i), \end{aligned}$$

because for every  $(i_1, \dots, i_r)$  in the sum  $\sum_{n,r}$ ,

$$E[h(X_{i_1}, \dots, X_{i_r}) | X_i] = \begin{cases} h_1(X_i) & \text{if } i \in \{i_1, \dots, i_r\} \\ E[h(X_{i_1}, \dots, X_{i_r})] = 0 & \text{otherwise,} \end{cases}$$

so that

$$\sum_{n,r} \mathbb{E}[h(X_{i_1}, \dots, X_{i_r})|X_i] = r(n-1)^{(r-1)} h_1(X_i) \text{ and}$$

$$\sqrt{n}(n-1)^{(r-1)}/n^{(r)} = n^{-1/2}.$$

Hence

- (i)  $T_n^* = rn^{-1/2} \sum_{i=1}^n h_1(X_i) \xrightarrow{\mathcal{L}} N(0, r^2 \xi_1)$ ,
- (ii)  $\text{Var}[T_n^*] = r^2 \xi_1$ ,
- (iii)  $\text{Var}[T_n] = n\text{Var}[U_n] = r^2 \xi_1 + o(1)$  as already shown, and
- (iv)  $\mathbb{E}\left[\left(T_n - T_n^*\right)^2\right] = \text{Var}[T_n] - \text{Var}[T_n^*] = o(1)$ ,

so that  $T_n = T_n^* + o_P(1)$ .

Thus  $T_n = \sqrt{n}(U_n - \theta) \xrightarrow{\mathcal{L}} N(0, r^2 \xi_1)$  by Slutsky's Theorem, completing the proof of [Theorem 8.1.3](#).

Similarly, in the two-sample problem, the Hájek projection of  $T_{m,n} = \sqrt{m+n}(U_{m,n} - \theta)$  is

$$\begin{aligned} T_{m,n}^* &= \sqrt{m+n} \left[ \frac{r(m-1)^{(r-1)}}{m^{(r)}} \sum_{i=1}^m h_{1,0}(X_i) + \frac{s(n-1)^{(s-1)}}{n^{(s)}} \sum_{j=1}^n h_{0,1}(Y_j) \right] \\ &= \sqrt{m+n} \left[ rm^{-1} \sum_{i=1}^m h_{1,0}(X_i) + sn^{-1} \sum_{j=1}^n h_{0,1}(Y_j) \right]. \end{aligned}$$

Now arguing as in one-sample problem, [Theorem 8.1.4](#) follows. □

### Examples

1. The Wilcoxon signed-rank statistic  $T_n$  can be written as

$$\begin{aligned} T_n &= \{n(n-1)\}^{-1} \sum_{i=1}^n S_i R_{n:i}^+ = \left(1/n^{(2)}\right) \left[ \sum_{i=1}^n S_i + \sum_{i \neq j=1}^n S_i I_{(0,\infty)}(|X_i| - |X_j|) \right] \\ &= \left(1/n^{(2)}\right) \sum_{n,2} (1/2) \left[ S_i I_{(0,\infty)}(|X_i| - |X_j|) + S_j I_{(0,\infty)}(|X_j| - |X_i|) \right] + O(1/n) \\ &= U_n + O(1/n), \text{ where} \\ U_n &= \left(1/n^{(2)}\right) \sum_{n,2} g(X_{i_1}, X_{i_2}), g(x_1, x_2) = I_{(0,\infty)}(x_1 + x_2) - 1/2. \end{aligned}$$

It is now easy to see that for this  $U$  statistic,

$$\theta(F) = \mathbb{E}_F[I_{(0,\infty)}(X_1 + X_2) - 1/2] = \int_{-\infty}^{\infty} [1 - F(-x)] dF(x) - 1/2$$

so under  $H_0$ :  $F(-x) = 1 - F(x)$  for all  $x$ ,  $\theta = 0$  and  $h(x_1, x_2) = g(x_1, x_2)$ . The variance of  $h_1(X_1) = E[h(X_1, X_2)|X_1] = (1/2) - \theta(F) - F(-X_1)$  is

$$\xi_1(F) = \int [1/2 - \theta(F) - F(-x)]^2 dF(x) > 0$$

iff  $F(-X)$  is not constant with probability 1. For continuous  $F$ , this holds iff  $0 < P_F[X < 0] < 1$ . Under this condition

$$\sqrt{n}[T_n/\{n(n-1)\} - \theta(F)] \xrightarrow{\mathcal{L}} N(0, 4\xi_1(F)).$$

In particular, let the common cdf of  $X_1, \dots, X_n$  be  $F(\cdot - \theta)$  where  $F$  has a pdf which is symmetric about 0. Then

$$\begin{aligned}\mu(\theta) &= \theta(F) = \int [1 - F(-x - \theta)] dF(x - \theta) - 1/2 \\ &= \int F(x + 2\theta) dF(x) - 1/2, \text{ and} \\ \sigma^2(\theta) &= 4\xi_1(F) = 4 \left[ \int F^2(x + 2\theta) dF(x) - \left\{ \int F(x + 2\theta) dF(x) \right\}^2 \right].\end{aligned}$$

2. The Mann-Whitney statistic  $U_{m,n} = (mn)^{-1} \sum_{i=1}^m \sum_{j=1}^n I_{(0,\infty)}(X_i - Y_j)$  based on independent rv's  $(X_1, \dots, X_m)$ ,  $(Y_1, \dots, Y_n)$ , where the  $X_i$ 's have a common cdf  $F$  and the  $Y_j$ 's have a common cdf  $G$  is a two-sample  $U$ -statistic with

$$\theta = \theta(F, G) = P_{F,G}[X_i < Y_j] = \int F dG = 1/2 \text{ if } F = G.$$

It is also easy to check that

$$h_{10}(x) = 1 - G(x) - E_F[1 - G(X)], \quad h_{01}(y) = F(y) - E_G[F(Y)],$$

so that  $\xi_{10} = \text{Var}_F[G(X)]$  and  $\xi_{01} = \text{Var}_G[F(Y)]$ . If  $F = G$ , then  $\xi_{10} = \xi_{01} = 1/12$ . Thus under  $H_0$ :  $F = G$ ,

$$\sqrt{m+n}(U_{m,n} - 1/2) \xrightarrow{\mathcal{L}} N(0, 1/(12\lambda(1-\lambda))),$$

if  $m, n \rightarrow \infty$  so that  $m/(m+n) \rightarrow \lambda \in (0, 1)$ . On the other hand, if  $G(x) = F(x - \theta)$ , then  $\mu(\theta) = \theta(F, G)$  and  $\sigma^2(\theta) = \lambda^{-1}\xi_{10} + (1-\lambda)^{-1}\xi_{01}$  are given by

$$\begin{aligned}\mu(\theta) &= \int F(x + \theta) dF(x) \text{ and} \\ \sigma^2(\theta) &= \lambda^{-1} \left[ \int F^2(x - \theta) dF(x) - \left\{ \int F(x - \theta) dF(x) \right\}^2 \right] \\ &\quad + (1-\lambda)^{-1} \left[ \int F^2(x + \theta) dF(x) - \left\{ \int F(x + \theta) dF(x) \right\}^2 \right].\end{aligned}$$

### 8.1.4 Asymptotic Comparison of Tests: Pitman's Approach

Let  $\theta$  be a numerical characteristic of a population (or a combination of populations), for which we want to test  $H_0: \theta = 0$  vs  $H_+: \theta > 0$ . Typically, a test statistic  $T_n$  based on a sample of size  $n$  would reject  $H_0$  in favor of  $H_+$  at level  $\alpha$  if  $T_n > c_n(\alpha)$ . Here we look into the problem of comparing the asymptotic powers of two sequences of tests based on  $\{T_n^{(1)}\}$  and  $\{T_n^{(2)}\}$  with critical values  $\{c_n^{(1)}(\alpha)\}$  and  $\{c_n^{(2)}(\alpha)\}$ , respectively, with same asymptotic Type I error probability at the same alternative. The following approach due to Pitman [See 69] attempts to make such a comparison by means of a measure of asymptotic relative efficiency (ARE) of one-test sequence with respect to the other. Since all reasonable tests are consistent (ie, having power at any  $\theta > 0$  tending to 1 as  $n \rightarrow \infty$  with Type I error probability fixed at  $0 < \alpha < 1$ ), we shall consider asymptotic powers of two-test sequences at alternatives  $\{\theta_n\}$  converging to the null hypothesis value  $\theta = 0$  at the rate of  $1/\sqrt{n}$ . The following theorem provides a formula for computing such asymptotic powers.

**Theorem 8.1.6** (Pitman). *Suppose that  $\{T_n\}$  and  $\{c_n\}$  are such that*

- (i)  $\lim_{n \rightarrow \infty} P_\theta[\sqrt{n}(T_n - \mu(\theta))/\sigma(\theta) \leq t] = \Phi(t)$  uniformly in a neighborhood of  $\theta = 0$ ,  $\mu(\cdot)$  is differentiable at 0 and  $\sigma(\cdot)$  is continuous at 0 with  $\sigma(0) > 0$  (here  $\Phi$  denotes the cdf of the standard normal),
- (ii)  $\lim_{n \rightarrow \infty} P_{\theta=0}[T_n \geq c_n] = \alpha$ .

Then

$$\lim_{n \rightarrow \infty} P_{\theta=\delta/\sqrt{n}}[T_n \geq c_n] = \Phi\left(\Phi^{-1}(\alpha) + \delta\mu'(0)/\sigma(0)\right).$$

*Proof.* For fixed  $\theta$ , the convergence in (i) is uniform in  $t$  by Polya's Theorem (Chapter 3, Theorem 3.2.5(VIII)), which will be used in the proof. Note that

$$\begin{aligned} \alpha &= \lim_{n \rightarrow \infty} P_0[\sqrt{n}(T_n - \mu(0))/\sigma(0) \geq \sqrt{n}(c_n - \mu(0))/\sigma(0)] \\ &= 1 - \Phi\left(\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(0))/\sigma(0)\right) = \Phi\left(-\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(0))/\sigma(0)\right). \end{aligned}$$

Hence

$$\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(0))/\sigma(0) = -\Phi^{-1}(\alpha).$$

Now

$$\begin{aligned} \lim_{n \rightarrow \infty} P_{\delta/\sqrt{n}}[T_n \geq c_n] &= \lim_{n \rightarrow \infty} P_{\delta/\sqrt{n}}[\sqrt{n}(T_n - \mu(\delta/\sqrt{n}))/\sigma(\delta/\sqrt{n}) \leq \sqrt{n}(c_n - \mu(\delta/\sqrt{n}))/\sigma(\delta/\sqrt{n})] \\ &= 1 - \Phi\left(\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(\delta/\sqrt{n}))/\sigma(\delta/\sqrt{n})\right) \\ &= \Phi\left(-\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(\delta/\sqrt{n}))/\sigma(\delta/\sqrt{n})\right). \end{aligned}$$

Finally,

$$\begin{aligned} -\lim_{n \rightarrow \infty} \sqrt{n}(c_n - \mu(\delta/\sqrt{n}))/\sigma(\delta/\sqrt{n}) &= -\lim_{n \rightarrow \infty} \{\sqrt{n}(c_n - \mu(0))/\sigma(0)\}\{\sigma(0)/\sigma(\delta/\sqrt{n})\} \\ &\quad + \lim_{n \rightarrow \infty} \{(\mu(\delta/\sqrt{n}) - \mu(0))/(\delta/\sqrt{n})\}\{\delta/\sigma(\delta/\sqrt{n})\} \\ &= \Phi^{-1}(\alpha) + \delta\mu'(0)/\sigma(0). \end{aligned}$$

□

In view of this theorem, if  $T_{n_i}^{(i)}$ ,  $i = 1, 2$  have the same asymptotic power at  $\delta_i/\sqrt{n_i}$ ,  $i = 1, 2$ , and if  $\delta_1/\sqrt{n_1} = \delta_2/\sqrt{n_2} = \theta_n$  (say), then

$$\lim_{n_i \rightarrow \infty} \beta_{T_{n_i}^{(i)}}(\theta_n) = \Phi\left(\Phi^{-1}(\alpha) + \delta_i \mu'_i(0)/\sigma_i(0)\right), \quad i = 1, 2,$$

must be equal, that is

$$\sqrt{n_2/n_1} = \delta_2/\delta_1 = \frac{\mu'_1(0)/\sigma_1(0)}{\mu'_2(0)/\sigma_2(0)}.$$

The ratio  $n_2/n_1$  is a measure of relative efficiency of  $\{T_n^{(1)}\}$  in comparison with  $\{T_n^{(2)}\}$ . For this reason, the ratio

$$e_{1,2} = \frac{\{\mu'_1(0)/\sigma_1(0)\}^2}{\{\mu'_2(0)/\sigma_2(0)\}^2}$$

is called the *asymptotic relative efficiency (ARE)* of  $\{T_n^{(1)}\}$  with respect to  $\{T_n^{(2)}\}$ . The quantity  $\{\mu'_i(0)/\sigma_i(0)\}^2$  is called the *asymptotic efficacy* of the sequence  $\{T_n^{(i)}\}$ .

**Example 8.1.2.** Let  $F$  be an unknown cdf with pdf  $f = F'$  which is symmetric about 0 and let  $G(x, \theta) = F(x - \theta)$ . Based on a random sample  $(X_1, \dots, X_n)$  from  $G(x, \theta)$ , we want to test  $H_0: \theta = 0$  vs  $H_+: \theta > 0$ . Consider three test statistics

$$T_n^{(1)} = \bar{X}_n/s_n, \quad T_n^{(2)} = \bar{S}_n, \quad \text{and} \quad T_n^{(3)} = \{n(n-1)\}^{-1} \sum_{i=1}^n S_i R_{n:i}^+,$$

for the  $t$ -test, the sign test, and the Wilcoxon singed-rank test. Then all three sequences  $\{T_n^{(j)}: j = 1, 2, 3\}$  satisfy the conditions of [Theorem 8.1.6](#) with

$$\mu_1(\theta) = \theta/\sigma(F), \quad \sigma_1^2(\theta) = 1 + \theta^2 h(\theta) \text{ where } h(\theta) = \{\mu_4(F) - \sigma^4(F)\}/\{4\sigma^6(F)\} \text{ for } T_n^{(1)},$$

$$\mu_2(\theta) = 1 - 2F(-\theta), \quad \sigma_2^2(\theta) = 4F(\theta)\{1 - F(\theta)\} \text{ for } T_n^{(2)}, \text{ and}$$

$$\mu_3(\theta) = \int F(x + 2\theta) dF(x) - 1/2, \quad \sigma_3^2(\theta) = 4\text{Var}_F[F(X + 2\theta)], \text{ where}$$

$$\sigma^2(F) = E_F[(X - \theta)^2] \text{ and } \mu_4(F) = E_F[(X - \theta)^4].$$

Thus the asymptotic efficacies of  $\{T_n^{(1)}\}$ ,  $\{T_n^{(2)}\}$ , and  $\{T_n^{(3)}\}$  are

$$\{\mu'_1(0)/\sigma_1(0)\}^2 = 1/\sigma^2(F), \quad \{\mu'_2(0)/\sigma_2(0)\}^2 = 4f^2(0), \text{ and}$$

$$\{\mu'_3(0)/\sigma_3(0)\}^2 = 12 \left\{ \int_{-\infty}^{\infty} f^2(x) dx \right\}^2.$$

If  $f$  is the pdf of  $N(0, 1)$ , then the AREs of  $\{T_n^{(2)}\}$  and  $\{T_n^{(3)}\}$  with respect to  $\{T_n^{(1)}\}$  are

$$e_{2,1} = 4f^2(0)\sigma^2(F) = 2/\pi \text{ and } e_{3,1} = 12 \left\{ \int_{-\infty}^{\infty} f^2(x) dx \right\}^2 \sigma^2(F) = 3/\pi.$$

## An Outline of Contiguity Theory

Let  $S_N$  be a test statistic based on observations  $X_1, \dots, X_N$ , iid as  $X$  distributed with pdf/pmf  $f(x; \theta)$  to test  $H_0: \theta = 0$  vs  $H_1: \theta > 0$ . In the next example, we shall illustrate the use of Contiguity Theory in deriving the asymptotic distribution of  $S_N$  under  $f(x; \delta/\sqrt{N})$  from its asymptotic distribution under  $f(x; 0)$  (see [4, 34]).

For  $N = 1, 2, \dots$ , let  $P_N$  and  $Q_N$  be probabilities on  $(\mathcal{X}_N, \mathcal{A}_N)$ . The sequence  $\{Q_N\}$  is said to be contiguous to  $\{P_N\}$  if for any sequence  $A_N \in \mathcal{A}_N$ ,  $\lim_{N \rightarrow \infty} P_N(A_N) = 0$  implies  $\lim_{N \rightarrow \infty} Q_N(A_N) = 0$ . Let  $L_N$  denote the likelihood ratio of  $Q_N$  to  $P_N$ .

For our purpose,  $\mathcal{X}_N$  is the sample space of a dataset  $(X_1, \dots, X_N)$ ,  $\mathcal{A}_N$  is a family of events in  $\mathcal{X}_N$ , and  $P_N, Q_N$  are joint distributions of  $(X_1, \dots, X_N)$  under two models such as a null hypothesis and a sequence of alternatives. For each  $N$ , let  $S_N$  be a test statistic based on  $(X_1, \dots, X_N)$  and suppose that we know the asymptotic distribution of  $\{S_N\}$  under  $\{P_N\}$ . The aim is to find the asymptotic distribution of  $\{S_N\}$  under  $\{Q_N\}$  from this by using  $\{L_N\}$  as a link between  $\{P_N\}$  and  $\{Q_N\}$ . Operationally, this is achieved by finding the asymptotic joint distribution of  $\{(S_N, \log L_N)\}$  under  $\{P_N\}$  when  $\{Q_N\}$  is contiguous to  $\{P_N\}$ .

The following results due to LeCam (LeCam's First and Third Lemma) provide the main tools:

- I. If  $\log L_N \xrightarrow{\mathcal{L}} N(-\sigma^2/2, \sigma^2)$  for some  $\sigma^2 > 0$  under  $\{P_N\}$ , then  $\{Q_N\}$  is contiguous to  $\{P_N\}$ .
- II. If  $\begin{pmatrix} S_N \\ \log L_N \end{pmatrix} \xrightarrow{\mathcal{L}} N_2 \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \right)$  with  $\mu_2 = -\sigma_2^2/2$  under  $\{P_N\}$ , then  
 $S_N \xrightarrow{\mathcal{L}} N(\mu_1 + \sigma_{12}, \sigma_1^2)$  under  $\{Q_N\}$ .

**Example 8.1.3.** Let  $X_1, \dots, X_N$  be iid as  $X$  distributed with pdf/pmf  $f(x; \theta)$  and let  $P_N, Q_N$  denote, respectively, the joint distribution of  $(X_1, \dots, X_N)$  under  $f(x; 0)$  and  $f(x; \delta/\sqrt{N})$ . The log-likelihood ratio of  $Q_N$  to  $P_N$  is

$$\log L_N = \sum_{i=1}^N \left[ \log f(X_i; \delta/\sqrt{N}) - \log f(X_i; 0) \right].$$

Under  $P_N$ , subject to regularity conditions,

$$\log L_N = \left( \delta/\sqrt{N} \right) \sum_{i=1}^N \dot{l}(X_i; 0) - \left( \delta^2/2 \right) I(f) + o_P(1), \text{ as } N \rightarrow \infty, \text{ where}$$

$$\dot{l}(x; 0) = \frac{\partial}{\partial \theta} \log f(x; \theta) \Big|_{\theta=0} \quad \text{with } E_{\theta=0} [\dot{l}(X; 0)] = 0, \text{ and}$$

$$\mathbb{E}_{\theta=0} \left[ l(X; 0)^2 \right] = I(f) \text{ (Fisher-information) [see Chapter 5, Section 5.2].}$$

Suppose that  $S_N^{(1)}$  and  $S_N^{(2)}$  are two statistics based on  $(X_1, \dots, X_N)$ , which under  $P_N$ , are asymptotically of the form

$$S_N^{(j)} = N^{-1/2} \sum_{i=1}^N \psi_j(X_i) + o_P(1), \quad j = 1, 2, \text{ as } N \rightarrow \infty, \text{ where}$$

$$\mathbb{E}_{\theta=0} \left[ \psi_j(X) \right] = 0, \quad \mathbb{E}_{\theta=0} \left[ \psi_j^2(X) \right] = \sigma_j^2, \quad \text{and} \quad \delta \mathbb{E}_{\theta=0} \left[ l(X; 0) \psi_j(X) \right] = \mu_j.$$

- (a) Find the asymptotic joint distribution of  $(S_N^{(j)}, \log L_N)$ ,  $j = 1, 2$ , under  $P_N$ . What do these asymptotic distributions imply?
- (b) Find the asymptotic distribution of  $S_N^{(j)}$  under  $Q_N$ ,  $j = 1, 2$ .
- (c) Find the Pitman ARE of  $S_N^{(1)}$  with respect to  $S_N^{(2)}$ .

*Solution.* We shall use the Contiguity Theory.

- (a) By the bivariate CLT and Slutsky's Theorem, for  $j = 1, 2$ , and denoting  $\mu_j = \delta \mathbb{E}_{\theta=0} \left[ l(X; 0) \psi_j(X) \right]$ ,

$$\begin{bmatrix} S_N^{(j)} \\ \log L_N \end{bmatrix} \xrightarrow{\mathcal{L}} N_2 \left( \begin{bmatrix} 0 \\ -(\delta^2/2)I(f) \end{bmatrix}, \begin{bmatrix} \sigma_j^2 & \mu_j \\ \mu_j & \delta^2 I(f) \end{bmatrix} \right), \text{ where}$$

$$\mu_j = \delta \mathbb{E}_{\theta=0} \left[ l(X_j; 0) \psi_j(X) \right].$$

- (b) Since  $\log L_N \xrightarrow{\mathcal{L}} N(-(\delta/2)I(f), \delta^2 I(f))$ ,  $\{Q_N\}$  is contiguous to  $\{P_N\}$  by LeCam's First Lemma, so LeCam's Third Lemma applies, by which

$$S_N^{(j)} \xrightarrow{\mathcal{L}} N(0 + \mu_j, \sigma_j^2) = N(\delta \mathbb{E}_{\theta=0} \left[ l(X; 0) \psi_j(X) \right], \sigma_j^2), \quad j = 1, 2.$$

- (c) The Pitman ARE of  $S_N^{(2)}$  with respect to  $S_N^{(1)}$  is

$$e_{2,1} = \frac{\mathbb{E}_0 \left[ l(X; 0) \psi_2(X) \right] / \sigma_2^2}{\mathbb{E}_0 \left[ l(X; 0) \psi_1(X) \right] / \sigma_1^2}.$$

## 8.2 Locally Most Powerful Rank Tests

Tests discussed in the last section were proposed on an ad hoc basis, on intuitive grounds, without any kind of optimality criterion in mind. Although some of these tests have good asymptotic properties under normality in terms of Pitman's ARE, their performances under other models will vary.

It should be mentioned that for data which one suspects to be normally distributed without being sure about it, Fisher and Yates [35] proposed the test statistic

$$\sum_{i=m+1}^{m+n} \mathbb{E}[\Phi^{-1}(U_{N:R_{N;i}})], \quad N = m + n$$

for the two-sample problem, where  $U_{N:1} < \dots < U_{N:N}$  are the order statistics in a random sample of size  $N$  from  $\text{Unif}(0, 1)$  and  $\Phi$  is the cdf of  $N(0, 1)$ . Further investigations on this and related problems were carried out by Hoeffding [36] and Terry [37].

The material presented in this section is mainly based on the development by Hájek and Šidák [4] aiming at construction of rank tests with a local optimality property for a general class of nonparametric hypotheses.

**Definition 8.2.1.** Let  $\mathcal{P}$  and  $\{q_\Delta, \Delta \geq 0\}$  be families of pdf's on  $\mathbb{R}^N$ , so that  $q_0 \in \mathcal{P}$  but  $\{q_\Delta, \Delta > 0\}$  is distinct from  $\mathcal{P}$ , and suppose that  $p$  is the joint pdf of rv's  $(X_1, \dots, X_N)$ . A rank test  $\Psi^*$  is a locally most powerful (LMP) rank test at level  $\alpha$  for  $H_0: p \in \mathcal{P}$  vs  $H_1: p = q_\Delta$ ,  $\Delta > 0$  based on  $(X_1, \dots, X_N)$  if

- (i)  $\mathbb{E}_p[\Psi^*(\mathbf{R}_N)] \leq \alpha$  for all  $P \in \mathcal{P}$  where  $\mathbf{R}_N$  is the vector of ranks of  $(X_1, \dots, X_N)$  and
- (ii) there exists  $\varepsilon > 0$  such that  $\mathbb{E}_p[\Psi^*(\mathbf{R}_N)] \geq \mathbb{E}_p[\Psi(\mathbf{R}_N)]$  for all  $\Psi$  satisfying (i) and for all  $p \in \{q_\Delta, 0 < \Delta \leq \varepsilon\}$ .

In this section, we shall construct LMP rank tests for several nonparametric hypothesis testing problems. These problems will involve the following families of probabilities

$$\begin{aligned} \mathcal{P}_0 &= \left\{ p: p(\mathbf{x}) = \prod_{i=1}^N f(x_i), f \text{ is an unknown pdf on } \mathbb{R} \right\}, \\ \mathcal{P}_1 &= \left\{ p: p(\mathbf{x}) = \prod_{i=1}^N f(x_i), f \text{ is an unknown symmetric pdf on } \mathbb{R} \right\}, \\ \mathcal{P}_2 &= \left\{ p: p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^N f(x_i)g(y_i), f \text{ and } g \text{ are unknown pdf's on } \mathbb{R} \right\}. \end{aligned}$$

The problems of testing  $p \in \mathcal{P}_0$  or  $\mathcal{P}_1$  or  $\mathcal{P}_2$  vs various alternatives  $\{q_\Delta, \Delta > 0\}$  will be taken up in the following.

### 8.2.1 Testing $H_0: p \in \mathcal{P}_0$ Against a General Alternative

*Regularity Conditions A.* Let  $\{f(x, \theta), \theta \in J\}$  be a family of pdf's on  $\mathbb{R}$  where  $J$  is an open interval in  $\mathbb{R}$  containing 0 and the following hold:

- (i)  $f(x, \theta)$  is absolutely continuous in  $\theta$  for almost all  $x$  (see Section 2.1 of Chapter I of Hájek and Šidák [4]),

(ii) for each  $\theta$  in a neighborhood of 0, the partial derivative

$$\dot{f}(x, \theta) = \lim_{h \rightarrow 0} \{f(x, \theta + h) - f(x, \theta)\}/h$$

exists for almost all  $x$ ,

$$(iii) \lim_{\theta \rightarrow 0} \int_{-\infty}^{\infty} |\dot{f}(x, \theta)| dx = \int_{-\infty}^{\infty} |\dot{f}(x, 0)| dx < \infty.$$

**Theorem 8.2.1.** Let  $p$  denote an unknown joint pdf of  $(X_1, \dots, X_N)$  and let  $q_{\Delta}(x_1, \dots, x_N) = \prod_{i=1}^N f(x_i, c_i \Delta)$ , where the family  $\{f(x, \theta), \theta \in J\}$  is specified and satisfies the regularity conditions A. Then the LMP rank test at level  $\alpha$  for  $H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q_{\Delta}, \Delta > 0$  is given by

$\Psi(\mathbf{r}) = 0, \text{ or } \gamma, \text{ or } 1$  according as

$$\sum_{i=1}^N c_i a_N(r_i, f) <, \text{ or } =, \text{ or } > k, \text{ where}$$

$$a_N(i, f) = E_{q_0} \left[ \frac{\dot{f}}{f}(X_{N:i}, 0) \right],$$

and  $0 \leq \gamma \leq 1$  and  $k$  are determined by  $E_p[\Psi(\mathbf{R})] = \alpha$  for all  $p \in \mathcal{P}_0$ .

Before going into the proof, we shall look into some special cases of the family  $\{f(x, \theta)\}$  which determines the score function  $a_N(i, f)$ .

### Location and Scale Families

**Definition 8.2.2.** The family  $\{f(x, \theta) = f(x - \theta), -\infty < \theta < \infty\}$  is called a location family and the family  $\{f(x, \theta) = e^{-\theta} f((x - \mu)e^{-\theta}), -\infty < \theta < \infty\}$  is called a scale family.

**Lemma 8.2.1.** Let  $f$  be an absolutely continuous pdf. Then

- (a) the location family satisfies Condition A if  $\int_{-\infty}^{\infty} |f'(x)| dx < \infty$ , and
- (b) the scale family satisfies Condition A if  $\int_{-\infty}^{\infty} |xf'(x)| dx < \infty$ .

In the location family with  $f(x, \theta) = f(x - \theta)$ , we have

$$\dot{f}(x, \theta) = -f'(x - \theta) \text{ and } \frac{\dot{f}}{f}(X_{n:i}, 0) = -\frac{f'}{f}(X_{N:i}).$$

Let  $F(x, \theta) = F(x - \theta)$  where  $F$  is the cdf corresponding to  $f$ . Then

$$F^{-1}(u, \theta) = \theta + F^{-1}(u), \quad 0 < u < 1, \text{ and}$$

$$\frac{\dot{f}}{f}(F^{-1}(u, 0), 0) = -\frac{f'}{f}(F^{-1}(u)) := \varphi(u, f).$$

Thus  $a_N(i, f) = E[\varphi(U_{N:i}, f)]$  where  $U_{N:1} < \dots < U_{N:N}$  are the order statistics in a random sample  $(U_1, \dots, U_N)$  from  $Unif(0, 1)$ .

In the scale family,

$$\begin{aligned}\dot{f}(x, \theta) &= e^{-\theta} f((x - \mu)e^{-\theta}) - (x - \mu)e^{-2\theta} f'((x - \mu)e^{-\theta}) \\ &= f(x, \theta) \left[ -1 - (x - \mu)e^{-\theta} \frac{f'}{f}((x - \mu)e^{-\theta}) \right], \text{ and} \\ F(x, \theta) &= \int_{-\infty}^x e^{-\theta} f((y - \mu)e^{-\theta}) dy = \int_{-\infty}^{(x-\mu)e^{-\theta}} f(y) dy = F((x - \mu)e^{-\theta}).\end{aligned}$$

Then  $F^{-1}(u, \theta) = \mu + e^\theta F^{-1}(u)$ ,  $0 < u < 1$ , and

$$\frac{\dot{f}}{f}(F^{-1}(u, 0), 0) = -1 - F^{-1}(u) \frac{f'}{f}(u) := \varphi_1(u, f).$$

The score function  $a_N(i, f)$  in this case is now denoted by

$$a_{1N}(i, f) = E[\varphi_1(U_{N:i}, f)].$$

The above results are summarized below.

Let  $F$  denote the cdf corresponding to the pdf  $f$  by means of which the location and scale families are defined and let  $U_{N:i}$  denote the  $i$ th-order statistic in a random sample of size  $N$  from  $\text{Unif}(0, 1)$ . Then the score functions in the location and the scale family are, respectively,

$$a_N(i, f) = E[\varphi(U_{N:i}, f)] \text{ where } \varphi(u, f) = -\frac{f'}{f}(F^{-1}(u)), \quad (3)$$

$$a_{1N}(i, f) = E[\varphi_1(U_{N:i}, f)] \text{ where } \varphi_1(u, f) = -1 - F^{-1}(u) \frac{f'}{f}(F^{-1}(u)). \quad (4)$$

Next consider the important case when

$$c_1 = \dots = c_m = 0 \text{ and } c_{m+1} = \dots = c_N = 1,$$

which corresponds to the *two-sample problem* in which  $(X_1, \dots, X_m)$  and  $(X_{m+1}, \dots, X_N)$  are independent random samples from two populations and we want to test whether they have the same pdf, or whether their pdf's differ in location (or scale). The test criterion for the LMP rank tests in the two-sample location problem and the two-sample scale problem are, respectively,

$$\sum_{i=m+1}^N a_N(r_i, f) \text{ and } \sum_{i=m+1}^N a_{1N}(r_i, f),$$

where the score functions are given by Eqs. (3) and (4).

We now give an outline of the proof of [Theorem 8.2.1](#). The proofs of three subsequent theorems dealing with LMP rank tests in which  $P_0$  is replaced by  $P_1$  or  $P_2$  in  $H_0$  and  $q_\Delta$  have other forms, will be concerned with the likelihood ratio of  $R_N$  or  $(R_n^+, S)$ , etc., namely,

$$P_{q_\Delta}[R_N = r]/P_{q_0}[R_N = r], \text{ etc.},$$

which will be written as

$$\begin{aligned} & 1 + N! [P_{q_\Delta}(\mathbf{R}_N = \mathbf{r}) - P_{q_0}(\mathbf{R}_N = \mathbf{r})] \\ &= 1 + N! \Delta \int_{\mathbf{R}_N = \mathbf{r}} \Delta^{-1} \left\{ \prod_{i=1}^N f(x_i, \Delta c_i) - \prod_{i=1}^N f(x_i, 0) \right\} d\mathbf{x} \end{aligned}$$

in the proof of [Theorem 8.2.1](#) and similarly in the other proofs. The integral in the expression above is an  $N$ -dimensional integral. The integrand is a difference of two products, which will be simplified by the identity

$$\prod_{i=1}^N a_i - \prod_{i=1}^N b_i = \sum_{k=1}^k (a_k - b_k) \prod_{j=1}^{k-1} a_j \prod_{j=k+1}^N b_j \quad (5)$$

resulting in a sum of  $k$  integrals. The crucial step is taking  $\lim_{\Delta \downarrow 0}$  under the integrals, where the regularity conditions come into play. The proof is then completed by some routine simplifications, using [Lemma 8.2.1](#).

*Proof of Theorem 8.2.1.* By N-P Lemma, it is enough to show that there exists  $\varepsilon > 0$  such that for all  $\Delta \in (0, \varepsilon]$ , the likelihood ratios

$$P_{q_\Delta}[\mathbf{R}_N = \mathbf{r}] / P_{q_0}[\mathbf{R}_N = \mathbf{r}] \quad \text{for all } \mathbf{r}$$

are in the same ascending order as the numbers  $\sum_{i=1}^N c_i a_N(r_i, f)$ . Since  $P_{q_0}[\mathbf{R}_N = \mathbf{r}] = 1/N!$  for all  $\mathbf{r}$ , the likelihood ratio equals

$$\begin{aligned} & 1 + N! \int_{\mathbf{R}_N = \mathbf{r}} \left\{ \prod_{i=1}^N f(x_i, \Delta c_i) - \prod_{i=1}^N f(x_i, 0) \right\} d\mathbf{x} = 1 + N! \Delta \sum_{k=1}^N c_k I_{Nk}(\mathbf{r}, f; \Delta), \text{ where} \\ & I_{Nk} = \int_{\mathbf{R}_N = \mathbf{r}} \left\{ \frac{f(x_k, \Delta c_k) - f(x_k, 0)}{\Delta c_k} \right\} \prod_{j=1}^{k-1} f(x_j, \Delta c_j) \prod_{j=k+1}^N f(x_j, 0) d\mathbf{x}, \end{aligned} \quad (6)$$

by Eq. (5). Hence the theorem will be proved by showing that

$$\begin{aligned} \lim_{\Delta \downarrow 0} I_{Nk}(\mathbf{r}, f; \Delta) &= a_N(r_k, f) = E_{q_0} \left[ \frac{\dot{f}}{f}(X_{N:k}, 0) | \mathbf{R}_N = \mathbf{r} \right] \\ &= E_{q_0} \left[ \frac{\dot{f}}{f}(X_{N:k}, 0) \right], \end{aligned}$$

using [Lemma 8.1.1](#), and because  $N! \prod_{j=1}^N f(x_j, 0)$  is the conditional pdf of  $\mathbf{X}$  given  $\mathbf{R}_N = \mathbf{r}$  under  $q_0$ .

The justification for taking  $\lim_{\Delta \downarrow 0}$  under the integral is provided by the regularity conditions on  $\{f(x, \theta), \theta \in J\}$ . We omit the technicalities of this demonstration, referring to Hájek and Šidák [4].  $\square$

### 8.2.2 One-Sample Location Problem, Assuming Symmetry

Here we shall consider rank-sign tests based on  $(\mathbf{R}_N^+, \mathbf{S})$ , the properties of which under symmetry are given in [Theorem 8.1.2](#), and a rank-sign test has been discussed in [Section 8.1.1](#).

**Theorem 8.2.2.** Suppose that  $(X_1, \dots, X_N)$  have joint pdf  $p(\mathbf{x})$  on  $R^N$ . Then the LMP rank-sign test at level  $\alpha$  for  $H_0: p \in \mathcal{P}_1$  vs  $H_1: p(\mathbf{x}) = q_\Delta(\mathbf{x}) = \prod_{i=1}^N f(x_i - \Delta)$ ,  $\Delta > 0$ , where  $f$  is a specified symmetric absolutely continuous pdf on  $R$  with  $\int_{-\infty}^{\infty} |f'(x)| dx < \infty$ , is given by

$$\Psi(\mathbf{r}, \mathbf{s}) = 0, \text{ or } \gamma, \text{ or } 1 \text{ according as}$$

$$\sum_{i=1}^N s_i a_N^+(r_i, f) <, \text{ or } = \text{ or } > k,$$

the constant  $k$  and  $0 \leq \gamma \leq 1$  being determined by  $E_{H_0}[\Psi(\mathbf{R}_N^+, \mathbf{S})] = \alpha$  and the score function is

$$\begin{aligned} a_N^+(i, f) &= E_{q_0}\left[-\frac{f'}{f}(|X|_{N:i})\right] = E_{q_0}\left[-\frac{f'}{f}\left(F^{-1}(1/2 + (1/2)U_{N:i})\right)\right] \\ &:= E_{q_0}[\varphi^+(U_{N:i}, f)]. \end{aligned}$$

*Proof.* We need to show that the likelihood ratios

$$(2^N N!) P_{q_\Delta}[(\mathbf{R}_N^+, \mathbf{S}) = (\mathbf{r}, \mathbf{s})]$$

for the  $2^N N!$  different  $(\mathbf{r}, \mathbf{s})$  are ordered in the same way as  $\sum_{i=1}^N s_i a_N^+(r_i, f)$  for all sufficiently small  $\Delta$ . As indicated earlier, we write these likelihood ratios as

$$1 + (2^N N!) \Delta \sum_{i=1}^N I_{Nk}(\mathbf{r}, \mathbf{s}, f; \Delta)$$

using Eq. (5) where the integrals  $I_{Nk}$  have the same form as in Eq. (6) except that the integration is over the set  $\{(\mathbf{R}_N^+, \mathbf{S}) = (\mathbf{r}, \mathbf{s})\}$ . By [Lemma 8.2.1\(a\)](#), the conditions on  $f$  allow taking  $\lim_{\Delta \downarrow 0}$  under these integrals, showing that

$$\lim_{\Delta \downarrow 0} I_{Nk}(\mathbf{r}, \mathbf{s}, f, \Delta) = \int_{(\mathbf{R}_N^+, \mathbf{S}) = (\mathbf{r}, \mathbf{s})} \left\{ -\frac{f'}{f}(x_k) \right\} f(\mathbf{x}) d\mathbf{x}.$$

Finally, since  $f$  is symmetric,

$$\frac{f'}{f}(x) = sign(x) \frac{f'}{f}(|x|).$$

Hence

$$\begin{aligned} \lim_{\Delta \downarrow 0} I_{Nk}(\mathbf{r}, \mathbf{s}, f, \Delta) &= s_k \int_{(\mathbf{R}_N^+, \mathbf{S}) = (\mathbf{r}, \mathbf{s})} \left\{ -\frac{f'}{f}(|x_k|) \right\} f(\mathbf{x}) d\mathbf{x} \\ &= \frac{s_k}{2^N N!} E_{q_0}\left[ -\frac{f'}{f}(|X_k|) \mid \mathbf{R}_N^+ = \mathbf{r}, \mathbf{S} = \mathbf{s} \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{s_k}{2^N N!} E_{q_0} \left[ -\frac{f'}{f} (|X_k|) | \mathbf{R}_N^+ = \mathbf{r} \right] \\
&= \frac{s_k}{2^N N!} E_{q_0} \left[ -\frac{f'}{f} (|X|_{N:i}) \right] \propto s_k a_N^+(r_k, f)
\end{aligned}$$

by [Lemma 8.1.1](#) and since  $|X_k|$  is independent of  $\mathbf{S}$  under  $q_0$ . This concludes the proof, leaving the verification of  $a_N^+(i, f) = E[\varphi^+(U_{N:i}, f)]$  as an exercise.  $\square$

### 8.2.3 Two-Sample Scale Problem, Assuming Symmetry

**Theorem 8.2.3.** *Let  $p(\mathbf{x})$  be the joint pdf of  $(X_1, \dots, X_N)$ . Then the LMP rank test for  $H_0: p \in \mathcal{P}_1$  vs  $H_1: p(\mathbf{x}) = q_\Delta(\mathbf{x}) = \prod_{i=1}^m f(x_i, 0) \prod_{i=m+1}^N f(x_i, \Delta)$ , where  $\{f(x, \theta) = e^{-\theta} f(xe^{-\theta})\}$  is a scale family with a specified symmetric absolutely continuous pdf  $f$  on  $\mathbb{R}$  having  $\int_{-\infty}^{\infty} |xf'(x)| dx < \infty$ , is given by*

$$\Psi(\mathbf{r}) = 0, \text{ or } \gamma, \text{ or } 1 \text{ according as}$$

$$\sum_{i=m+1}^N a_{1N}^+(r_i, f) <, \text{ or } =, \text{ or } > k,$$

the constants  $k$  and  $0 \leq \gamma \leq 1$  being determined by  $E_{H_0}[\Psi(\mathbf{R}_N^+)] = \alpha$  and the score function is

$$\begin{aligned}
a_{1N}^+(i, f) &= E_{q_0} \left[ -1 - |X|_{N:i} \frac{f'}{f} (|X|_{N:i}) \right] \\
&= E \left[ -1 - F^{-1}(1/2 + 1/2 U_{N:i}) \frac{f'}{f} \left( F^{-1}(1/2 + 1/2 U_{N:i}) \right) \right] \\
&:= E[\varphi_1^+(U_{N:i}, f)].
\end{aligned}$$

*Proof.* This is a special case of [Theorem 8.2.1](#) with  $c_1 = \dots = c_m = 0$ ,  $c_{m+1} = \dots = c_N = 1$ , and  $f(x, \theta) = e^{-\theta} f(xe^{-\theta})$ . [Lemma 8.2.1\(b\)](#) and the conditions on  $f$  imply Condition A for the validity of [Theorem 8.2.1](#). We have already seen that for this  $f(x, \theta)$ ,

$$\begin{aligned}
\dot{f}(x, \theta)/f(x, \theta) &= -1 - xe^{-\theta} \frac{f'}{f}(xe^{-\theta}), \text{ so} \\
\dot{f}(x, 0) &= -1 - x \frac{f'}{f}(x).
\end{aligned}$$

Also, since  $f$  is symmetric,

$$x = sign(x)|x| \text{ and } \frac{f'}{f}(x) = sign(x) \frac{f'}{f}(|x|).$$

Hence the score function is

$$a_{1N}^+(i, f) = E_{q_0} \left[ -1 - |X|_{N:i} \frac{f'}{f} (|X|_{N:i}) \right].$$

The verification of  $a_{1N}^+(i, f) = E[\varphi_1^+(U_{N:i}, f)]$  is left as an exercise.  $\square$

### 8.2.4 Test for Independence in a Bivariate Population

**Theorem 8.2.4.** Let  $p(\mathbf{x}, \mathbf{y})$  be the joint pdf of  $\{(X_i, Y_i): i = 1, \dots, N\}$  and  $\mathbf{R}_N, \mathbf{R}'_N$  the vectors of ranks of  $(X_1, \dots, X_N)$  and  $(Y_1, \dots, Y_N)$ , respectively. Then the LMP rank test at level  $\alpha$  based on  $\mathbf{R}_N$  and  $\mathbf{R}'_N$  for  $H_0: p \in \mathcal{P}_2$  vs  $H_1: p(\mathbf{x}, \mathbf{y}) = q_\Delta(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^N h_\Delta(x_i, y_i)$ ,  $\Delta > 0$ , where  $h_\Delta(x, y) = \int f(x - \Delta z)g(y - \Delta z) dM(z)$  with specified pdf's  $f$  and  $g$  on  $\mathbb{R}$  having continuous  $f', g'$  and specified cdf  $M$  on  $\mathbb{R}$  with finite variance  $\sigma_M^2$  is

$$\Psi(\mathbf{r}, \mathbf{r}') = 0, \text{ or } \gamma, \text{ or } 1 \text{ according as}$$

$$\sum_{i=1}^N a_N(r_i, f)a_N(r'_i, g) <, \text{ or } =, \text{ or } > k,$$

the constants  $k$  and  $0 \leq \gamma \leq 1$  being determined by  $E_{H_0}[\Psi(\mathbf{R}_N, \mathbf{R}'_n)] = \alpha$  and the score functions  $a_N(i, f), a_N(i, g)$  are as in Eq. (3).

*Remark 8.2.1.* Under  $q_\Delta$ , the  $(X_i, Y_i)$ 's are iid as  $(X, Y) = (X^* + \Delta Z, Y^* + \Delta Z)$  where  $X^*$  with pdf  $f$ ,  $Y^*$  with pdf  $g$ , and  $Z$  with cdf  $M$  are mutually independent.

In the proof of this theorem, we shall use

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= \prod_{i=1}^N \{f_\Delta(x_i)g_\Delta(y_i)\} \text{ with} \\ f_\Delta(x) &= \int f(x - \Delta z) dM(z) \text{ and } g_\Delta(y) = \int g(y - \Delta z) dM(z). \end{aligned}$$

Note that this  $p(\mathbf{x}, \mathbf{y}) \in \mathcal{P}_2$  and call this distribution  $Q_{0,\Delta}$ .

**Lemma 8.2.2.**  $\lim_{\Delta \downarrow 0} \Delta^{-2} [h_\Delta(x, y) - f_\Delta(x)g_\Delta(y)] = f'(x)g'(y)\sigma_M^2$ .

*Proof.* By algebraic rearrangements, we can write

$$\begin{aligned} h_\Delta(x, y) - f_\Delta(x)g_\Delta(y) &= \iint [f(x - \Delta z)g(y - \Delta z) - f(x - \Delta z)g(y - \Delta z')] dM(z) dM(z') \\ &= (1/2) \iint \{f(x - \Delta z) - f(x - \Delta z')\}\{g(y - \Delta z) - g(y - \Delta z')\} dM(z) dM(z'). \end{aligned}$$

Hence

$$\begin{aligned} \Delta^{-2} [h_\Delta(x, y) - f_\Delta(x)g_\Delta(y)] &= \iint \left\{ \frac{f(x - \Delta z) - f(x - \Delta z')}{\Delta(z - z')} \right\} \left\{ \frac{g(y - \Delta z) - g(y - \Delta z')}{\Delta(z - z')} \right\} (1/2)(z - z')^2 dM(z) dM(z'). \end{aligned}$$

Let  $A(\Delta, \delta)$  be the part of this integral over  $\{|z|, |z'| \leq \delta/\Delta\}$  and let  $R(\Delta, \delta)$  be the remainder. Then  $A(\Delta, \delta)$  can be made arbitrarily close to

$$f'(x)g'(y)(1/2) \int_{|z| \leq \delta/\Delta} \int_{|z'| \leq \delta/\Delta} (z - z')^2 dM(z) dM(z')$$

by making  $\delta > 0$  sufficiently small, and then this tends to  $f'(x)g'(y)\sigma_M^2$  as  $\Delta \rightarrow 0$ . Finally, with  $\|f\|, \|g\|$  as sup norms of  $f, g$  and  $C = \{|z|, |z'| \leq \delta/\Delta\}^c$ ,

$$\begin{aligned} R(\Delta, \delta) &\leq \left[ 2\|f\| \|g\|/\Delta^2 \right] \iint_C dM(z) dM(z') \\ &\leq \left[ 4\|f\| \|g\|/\delta^2 \right] (\delta/\Delta)^2 P_M[|Z| > \delta/\Delta] \\ &\leq \left[ 4\|f\| \|g\|/\delta^2 \right] \int_{|z|>\delta/\Delta} z^2 dM(z) \rightarrow 0 \end{aligned}$$

for arbitrary  $\delta$  as  $\Delta \rightarrow 0$ . □

*Proof.* As in the proofs of [Theorems 8.2.1–8.2.3](#), we consider the likelihood ratio

$$\begin{aligned} L(\mathbf{r}, \mathbf{r}') &= \frac{Q_\Delta(\mathbf{R} = \mathbf{r}, \mathbf{R}' = \mathbf{r}')}{Q_{0\Delta}(\mathbf{R} = \mathbf{r}, \mathbf{R}' = \mathbf{r}')} \\ &= 1 + (N!)^2 \Delta^2 \int_{\mathbf{R}=\mathbf{r}} \int_{\mathbf{R}'=\mathbf{r}'} \Delta^{-2} \left[ \prod_{i=1}^N h_\Delta(x_i, y_i) - \prod_{i=1}^N \{f_\Delta(x_i)g_\Delta(y_i)\} \right] d\mathbf{x} d\mathbf{y} \\ &= 1 + N! \Delta^2 \sum_{k=1}^N I_{Nk}(\mathbf{r}, \mathbf{r}', f, g; \Delta) \end{aligned}$$

using Eq. (5), where

$$I_{Nk} = \int_{\mathbf{R}=\mathbf{r}} \int_{\mathbf{R}'=\mathbf{r}'} \frac{h_\Delta(x_k, y_k) - f_\Delta(x_k)g_\Delta(y_k)}{\Delta^2} \prod_{j=1}^{k-1} h_\Delta(x_j, y_j) \prod_{j=k+1}^N \{f_\Delta(x_j)g_\Delta(y_j)\} d\mathbf{x} d\mathbf{y}.$$

Since the likelihood ratios  $L(\mathbf{r}, \mathbf{r}')$  are ordered as  $\sum_{k=1}^N I_{Nk}(\mathbf{r}, \mathbf{r}', f, g; \Delta)$ , the theorem will be proved by showing that for each  $(N, K)$ ,

$$(N!)^2 \lim_{\Delta \downarrow 0} I_{Nk}(\mathbf{r}, \mathbf{r}', f, g; \Delta) = a_N(r_i, f) a_N(r'_i, g).$$

By [Lemma 8.2.2](#), for each  $(N, k)$  the integrand of  $I_{Nk}(\mathbf{r}, \mathbf{r}', f, g; \Delta)$  converges to (as  $\Delta \rightarrow 0$ )

$$\frac{f'(x_k)}{f}(x_k) \frac{g'(y_k)}{g}(y_k) \prod_{j=1}^N \{f(x_j)g(y_j)\}.$$

Now taking  $\lim_{\Delta \downarrow 0}$  under the integral in  $I_{Nk}$ , for the justification of which we refer to Hájek and Šidák [4], we have

$$\begin{aligned}
& (N!)^2 \lim_{\Delta \downarrow 0} I_{Nk}(\mathbf{r}, \mathbf{r}', f, g; \Delta) \\
&= \sigma_M^2 \int_{\mathbf{R}=\mathbf{r}} \int_{\mathbf{R}'=\mathbf{r}'} \frac{f'}{f}(x_k) \frac{g'}{g}(y_k) \frac{\prod_{j=1}^N \{f(x_j)g(y_j)\}}{(1/N!)^2} dx dy \\
&= \sigma_M^2 E_{q_0} \left[ \frac{f'}{f}(X_k) \frac{g'}{g}(Y_k) | \mathbf{R}_N = \mathbf{r}, \mathbf{R}'_N = \mathbf{r}' \right] \\
&= \sigma_M^2 E_{q_0} \left[ \frac{f'}{f}(X_k) | \mathbf{R}_N = \mathbf{r} \right] E_{q_0} \left[ \frac{g'}{g}(Y_k) | \mathbf{R}'_N = \mathbf{r}' \right] \\
&= \sigma_M^2 E_{q_0} \left[ \frac{f'}{f}(X_{N:r_k}) \right] E_{q_0} \left[ \frac{g'}{g}(Y_{N:r'_k}) \right] = \sigma_M^2 a_N(r_k, f) a_N(r'_k, g),
\end{aligned}$$

using Lemma 8.1.1. □

### 8.2.5 Specific Rank Tests Using Approximate Scores

First, note that the scores  $a_N(i, f), a_{1N}(i, f), \dots$  can only be obtained from tables (if available) due to their complicated expressions. For example, there are tables for  $E[\Phi^{-1}(U_{N:i})]$ . However, since  $E[U_{N:i}] = i/(N+1)$  and  $\text{Var}[U_{N:i}] = i(N-i+1)/\{(N+1)^2(N+2)\} \rightarrow 0$  as  $N \rightarrow \infty$ , the distribution of  $U_{N:i}$  is concentrated near  $i/(N+1)$ ; so

$$a_N(i, f) = E[\varphi(U_{N:i}, f)] \approx \varphi(E[U_{N:i}], f) = \varphi(i/(N+1), f)$$

if  $\varphi$  is sufficiently smooth near  $u = i/(N+1)$ .

We now look at a number of specific problems.

#### I. Two-sample location

- (a) Normal:  $\varphi(u, f) = \Phi^{-1}(u)$ ,  $a_N(i, f) = E[\Phi^{-1}(U_{N:i})] \approx \Phi^{-1}(i/(N+1))$ .
- (b) Logistic:  $\varphi(u, f) = 2u - 1$ ,  $a_N(i, f) = E[2U_{N:i} - 1] = (2/(N+1))i - 1$ , or equivalently,  $i$  (Wilcoxon test).
- (c) Double exponential:  $\varphi(u, f) = \text{sign}(2u - 1)$ ,

$$a_N(i, f) = E[\text{sign}(2U_{N:i} - 1)] \approx \text{sign}(2i/(N+1) - 1) = \text{sign}(i - (N+1)/2).$$

The test statistic  $\sum_{i=m+1}^N \text{sign}(r_i - (N+1)/2)$ , or equivalently,

$$\sum_{i=m+1}^N (1/2)[\text{sign}(r_i - (N+1)/2) + 1]$$

# observations in the second sample exceeding the pooled median (Median test).

#### II. Two-sample scale

- (a) Normal:  $1 + a_{1N}(i, f) = E[\{\Phi^{-1}(U_{N:i})\}^2] \approx \{\Phi^{-1}(i/(N+1))\}^2$ .
- (b) Cauchy-type tail:  $f(x) = \frac{1}{2(1+|x|)^2}$ ;  $1 + \varphi_1(u, f) = 2|2u - 1|$ ,

$$1 + a_{1N}(i, f) = 2E[|2U_{N:i} - 1|]$$

has led to several statistics.

#### III. One-sample location, assuming symmetry

- (a) Normal:  $a_N^+(i, f) = E[\Phi^{-1}(1/2 + (1/2)U_{N:i})] \approx \Phi^{-1}(1/2 + (1/2)i/(N+1))$ .

- (b) Logistic:  $\varphi(u, f) = 2u - 1$ ,  $\varphi^+(u, f) = \varphi(1/2 + (1/2)u, f) = u$ ,  
 $a_N^+(i, f) = E[U_{N:i}] = i/(N+1)$ . Test statistic is  $\sum_{X_i > 0} R_{N:i}^+/(N+1)$ , or equivalently,  
 $\sum_{X_i > 0} R_{N:i}^+$  (Wilcoxon signed-rank test).
- IV. Bivariate independence: For  $f$  and  $g$  both logistic, test statistic is Spearman's rank correlation  $\rho_S$ .

The verification of these scores and resulting test statistics are left as exercises.

### 8.2.6 Asymptotic Distribution of Test Statistics for LMP Rank Tests

We need the distribution under  $H_0$  of the test statistic of an LMP rank test in order to find the critical value at a given level  $\alpha$ . In some cases, these statistics have special forms (at least approximately), for which exact distributions may be manageable (see [Section 8.1.2](#)), while in some other cases, asymptotic distributions can be derived by the  $U$ -statistic approach (see [Section 8.1.3](#)).

We first state without proof the following theorem from Hájek and Šidák [4], providing the asymptotic null distributions of a large class of rank statistics such as those constructed in [Theorem 8.2.1](#).

**Theorem 8.2.5.** Suppose that  $\varphi$  is a square-integrable function on  $[0, 1]$  with  $\bar{\varphi} = \int_0^1 \varphi(u) du$  and that the sequence  $\{c_N\}$  satisfies

$$\lim_{N \rightarrow \infty} \frac{\sum_{i=1}^N (c_{Ni} - \bar{c}_N)^2}{\max_{1 \leq i \leq N} (c_{Ni} - \bar{c}_N)^2} = \infty, \text{ where } \bar{c}_N = N^{-1} \sum_{i=1}^N c_{Ni}.$$

Let  $S_{c_N} = \sum_{i=1}^N c_{Ni} a_N^\varphi(R_{N:i})$ , where  $\mathbf{R}_N$  is the rank vector of  $(X_1, \dots, X_N)$  having joint pdf  $p(x) \in \mathcal{P}_0$  and  $a_N^\varphi(i) = E[\varphi(U_{N:i})]$ . Then  $(S_{c_N} - \mu_{c_N})/\sigma_{c_N} \xrightarrow{\mathcal{L}} N(0, 1)$ , where

$$\mu_{c_N} = c_N \sum_{i=1}^N a_N^\varphi(i) \text{ and } \sigma_{c_N}^2 = \sum_{i=1}^N (c_{Ni} - \bar{c}_N)^2 \int_0^1 \{\varphi(u) - \bar{\varphi}\}^2 du,$$

which is assumed to be positive.

In the two-sample case with  $c_{N1} = \dots = c_{Nm_N} = 0$ ,  $c_{N,m_N+1} = \dots = c_{N,m_N+n_N} = c_{NN} = 1$ ,  $\bar{c}_N = n_N/N = 1 - \lambda_N$ , the condition on  $\{c_N\}$  holds if  $m_N/N \rightarrow \lambda \in (0, 1)$ . Here  $S_{c_N} = \sum_{i=m_N+1}^N a_N^\varphi(R_{N:i})$ ,  $\mu_{c_N} = (1 - \lambda_N) \sum_{i=1}^N a_N^\varphi(i)$ , and  $\sigma_{c_N}^2 = N\lambda_N(1 - \lambda_N) \int_0^1 \{\varphi(u) - \bar{\varphi}\}^2 du$ .

We now discuss a different approach by Chernoff and Savage [38] for deriving the asymptotic distribution of the LMP rank statistic in the two-sample location problem under the null hypothesis as well as under location-shift.

Let  $(X_1, \dots, X_m)$  and  $(Y_1, \dots, Y_n)$  be independent random samples from populations with cdf's  $F$  and  $G$ , respectively, on  $\mathbb{R}$  and let the empirical cdf's of the two samples be denoted by

$$F_m(x) = m^{-1} \sum_{i=1}^m I_{(-\infty, x]}(X_i) \text{ and } G_n(y) = n^{-1} \sum_{i=1}^n I_{(-\infty, y]}(Y_i).$$

Next let

$$H(x) = \lambda_N F(x) + (1 - \lambda_N) G(x) \text{ and } H_N(x) = \lambda_N F_m(x) + (1 - \lambda_N) G_n(x),$$

where  $N = m + n$  and  $m/N = \lambda_N$ . Finally, let  $\{J_N\}$  be a sequence of functions on  $[0, 1]$  by suitably extending  $J_N(i/N) = a_N(i)$ ,  $i = 1, \dots, N$  (eg, as a step function), which converges to a function  $J$  obeying certain conditions to be stated later.

The main idea of the Chernoff-Savage approach is to recognize that

$H_N(X_i) = \text{empirical cdf of the combined sample evaluated at } X_i = R_{N:i}/N$  is the rank statistic of interest,

where  $R_{N:i} = \text{rank of } X_i \text{ in the combined sample}$ , so if  $J_N(i/N) = a_N(i)$ , then

$$T_N = \int J_N[H_N(x)] dF_m(x) = m^{-1} \sum_{i=1}^m a_N(R_{N:i}).$$

To find the asymptotic distribution of  $T_N$ , write

$$\begin{aligned} T_N - \int J[H(x)] dF(x) \\ = \int \{J[H_N(x)] - J[H(x)]\} dF_m(x) + \int J[H(x)] d[F_m(x) - F(x)] \\ + \int \{J_N[H_N(x)] - J[H_N(x)]\} dF_m(x) \end{aligned} \quad (7)$$

On the right-hand side of Eq. (7), expand the integrand of the first term as

$$\{H_N(x) - H(x)\}J'[H(x)] + (1/2)\{H_N(x) - H(x)\}^2 J''[\alpha_N H_N(x) + (1 - \alpha_N)H(x)],$$

$0 \leq \alpha_N \leq 1$  and make rearrangement of all the terms. Thus

$$\begin{aligned} T_N - \int J[H] dF &= \int [\lambda_N(F_m - F) + (1 - \lambda_N)(G_n - G)]J'(H) dF \\ &\quad + \int J(H)d(F_m - F) + R_{N1} + R_{N2} + R_{N3}, \end{aligned} \quad (8)$$

where

$$\begin{aligned} R_{N1} &= \int (H_N - H)J'(H)d(F_m - F), \\ R_{N2} &= (1/2) \int (H_N - H)^2 J'(\alpha_N H_N + (1 - \alpha_N)H) dF_m, \text{ and} \\ R_{N3} &= \int \{J_N(H_N) - J(H_N)\} dF_m. \end{aligned}$$

Integrating by parts, the main terms of Eq. (8) become

$$\begin{aligned} &\int [\lambda_N(F_m - F) + (1 - \lambda_N)G_n - G]J'(H) dF \\ &\quad - \int (F_m - F)J'(H)d[\lambda_N F + (1 - \lambda_N)G] \end{aligned}$$

$$= -(1 - \lambda_N) \int (F_m - F) J'(H) dG + (1 - \lambda_N) \int (G_n - G) J'(H) dF.$$

Now letting

$$J'[H(x)] dG(x) = dB(x), J'[H(x)] dF(x) = dB^*(x),$$

and integrating by parts, we are led to

$$\begin{aligned} T_N - \int J[H(x)] dF(x) \\ &= -(1 - \lambda_N) \int [F_m(x) - F(x)] dB(x) + (1 - \lambda_N) \int [G_n(x) - G(x)] dB^*(x) + R_N \\ &= (1 - \lambda_N) \int B(x) d[F_m(x) - F(x)] - (1 - \lambda_N) \int B^*(x) d[G_n(x) - G(x)] + R_N \\ &= (1 - \lambda_N) \left[ m^{-1} \sum_{i=1}^m \{B(X_i) - EB(X_i)\} - n^{-1} \sum_{i=1}^n \{B^*(Y_i) - EB^*(Y_i)\} \right] + R_N. \end{aligned}$$

The remainder term  $R_N = R_{N1} + R_{N2} + R_{N3}$  can be shown to be  $o_P(N^{-1/2})$  under the following conditions:

1.  $J(u) = \lim_{N \rightarrow \infty} J_N(u)$  exists for  $0 < u < 1$  and is not a constant.
2.  $\int_{\{x: 0 < H_N(x) < 1\}} \{J_N[H_N(x)] - J[H_N(x)]\} dF_m(x) = o_P(N^{-1/2})$ .
3.  $J_N(1) = O(N^{-1/2})$ .
4.  $|J^{(i)}(u)| = \left| \frac{d^i J}{du^i} \right| \leq \text{constant} |u(1-u)|^{-i-1/2+\delta}$ ,  $i = 0, 1, 2$  for some  $\delta > 0$ .

Also, since  $H = \lambda_N F + (1 - \lambda_N) G$  depends on  $N$ , so do the iid sequences  $\{B(X_i), i = 1, \dots, m\}$  and  $\{B^*(Y_i), i = 1, \dots, n\}$ . The asymptotic normality of  $\sqrt{N}[T_N - \int J[H(x)] dF(x)]$  would, therefore, have to be justified for this triangular array situation. This justification is also provided by the above conditions.

We now arrive at the following theorem.

**Theorem 8.2.6** (Chernoff-Savage). *Under Conditions 1–4, if  $0 < \lambda_0 \leq \lambda_N \leq 1 - \lambda_0 < 1$  for all  $N$ , then*

$$\frac{\sqrt{N} \left[ T_N - \int_{-\infty}^{\infty} J[H(x)] dF(x) \right]}{\sqrt{(1 - \lambda_N) \left\{ \left( \frac{1 - \lambda_N}{\lambda_N} \right) \text{Var}[B(X_1)] + \text{Var}[B^*(Y_1)] \right\}}} \xrightarrow{\mathcal{L}} N(0, 1).$$

To make the statement of the above theorem explicit, we need  $\text{Var}[B(X_1)]$  and  $\text{Var}[B^*(Y_1)]$ . Let  $F_1(x) = I(X_1 \leq x)$  and  $G_1(y) = I(Y_1 \leq y)$ . Then

$$E[F_1(x)] = F(x), \text{Cov}[F_1(x), F_1(y)] = \min\{F(x), F(y)\} - F(x)F(y)$$

and likewise for  $E[G_1(x)]$  and  $\text{Cov}[G_1(x), G_1(y)]$ . Using these, we have

$$\text{Var}[B(X_1)] = 2 \iint_{-\infty < x < y < \infty} F(x)\{1 - F(y)\} J'[H(x)] J'[H(y)] dG(x) dG(y) \quad (9)$$

and  $\text{Var}[B^*(Y_1)]$  is obtained by interchanging  $F$  and  $G$  in the above formula for  $\text{Var}[B(X_1)]$ . The asymptotic variance

$$V_N = (1 - \lambda_N) \left\{ \left( \frac{1 - \lambda_N}{\lambda_N} \right) \text{Var}[B(X_1)] + \text{Var}[B^*(Y_1)] \right\}$$

can now be written in an explicit form.

In the two-sample problem, with  $G(x) = F(x - \theta_N)$  and  $F' = f$ , the asymptotic normality holds uniformly in  $0 < \lambda_0 \leq \lambda_N \leq 1 - \lambda_0 < 1$  and  $\theta_N$  in some neighborhood of 0. If  $\theta_N \rightarrow 0$ , then

$$\begin{aligned} \left( \frac{\lambda_N}{1 - \lambda_N} \right) V_N &\rightarrow 2 \iint_{x < y} F(x) \{1 - F(y)\} J'[F(x)] J'[F(y)] dF(x) dF(y) \\ &= 2 \iint_{0 < u < v < 1} u(1 - v) J'(u) J'(v) du dv \\ &= 2 \iiint_{0 < s < u < v < t < 1} J'(u) J'(v) du dv ds dt \\ &= \iint_{0 < s < t < 1} [J^2(s) + J^2(t) - 2J(s)J(t)] ds dt \\ &= \int_0^1 J^2(t) dt - \left\{ \int_0^1 J(t) dt \right\}^2. \end{aligned} \tag{10}$$

Verifying the details of the derivations of Eqs. (9) and (10) are left as exercises.

Going back to the original problem with

$$J_N(i/N) = a_N(i) = \text{E}\left[K^{-1}(U_{N:i})\right],$$

$K$  being a strictly increasing absolutely continuous cdf, we need  $J = K^{-1}$  to satisfy Condition 4 of Theorem 8.2.6 and  $\lambda_N \rightarrow \lambda \in (0, 1)$ . Then  $\sqrt{N}[T_N - \mu(\theta)]/\sigma(\theta) \xrightarrow{\mathcal{L}} N(0, 1)$ , where  $\mu(\theta)$  and  $\sigma^2(\theta)$  are obtained by letting  $J = K^{-1}$  in  $\int J[H(x)] dF(x)$  and the formula for  $V_N$  using Eq. (9) and its counterpart for  $\text{Var}[B^*(Y_1)]$ . From these, we get

$$\begin{aligned} \mu'(0) &= -(1 - \lambda) \int_{-\infty}^{\infty} (K^{-1})' [F(x)] f(x) dF(x), \text{ and} \\ \sigma^2(0) &= \left( \frac{1 - \lambda}{\lambda} \right) \left[ \int_0^1 \left\{ K^{-1}(u) \right\}^2 du - \left\{ \int_0^1 K^{-1}(u) du \right\}^2 \right] = \left( \frac{1 - \lambda}{\lambda} \right) \sigma_K^2, \end{aligned}$$

if  $\sigma_K^2 = \text{Var}_K[X] < \infty$ . Hence the Pitman asymptotic efficacy of  $T_N$  is

$$\{\mu'(0)\}^2 / \sigma^2(0) = \frac{\lambda(1 - \lambda)}{\sigma_K^2} \left\{ \int_{-\infty}^{\infty} K^{-1'}[F(x)] f(x) dF(x) \right\}^2.$$

Now take  $J = K^{-1} = \Phi^{-1}$ . If  $F(x) = \Phi((x - a)/b)$  and we take  $K = \Phi$  in  $T_N$ , then

$$\begin{aligned}\int_{-\infty}^{\infty} K^{-1'}[F(x)]f(x) dF(x) &= b^{-1} \int_{-\infty}^{\infty} \Phi^{-1'}[\Phi(t)]\phi(t) d\Phi(t) \\ &= b^{-1} \int_{-\infty}^{\infty} d\Phi(t) = b^{-1},\end{aligned}$$

$$\text{so } \{\mu'(0)\}^2/\sigma^2(0) = \lambda(1 - \lambda)/b^2.$$

For  $F(x) = \Phi((x - a)/b)$ , the two-sample  $t$ -test has the same asymptotic efficacy. Thus the Pitman ARE of the LMP rank test with score function  $\Phi^{-1}$  with respect to the  $t$ -test = 1.

## 8.3 Tests Based on Empirical Distribution Function

Let  $X_1, \dots, X_n, \dots$  be a sequence of iid rv's with common cdf  $F$  which we assume to be continuous and strictly increasing. The random function

$$F_n(x) = n^{-1} \sum_{i=1}^n I_{(-\infty, x]}(X_i), \quad -\infty < x < \infty,$$

called the empirical distribution function (edf) based on  $X_1, \dots, X_n$  has already been defined in [Chapter 3](#).

**Theorem 8.3.1** (Glivenko-Cantelli).  $\sup_x |F_n(x) - F(x)| \rightarrow 0$  a.s., as  $n \rightarrow \infty$ .

*Proof.*

- (i) If  $|F_n(a) - F(a)| < \varepsilon/2$ ,  $|F_n(b) - F(b)| < \varepsilon/2$  and  $F(b) - F(a) < \varepsilon/2$  for  $a < b$ , then  $|F_n(x) - F(x)| < \varepsilon$  for all  $a \leq x \leq b$  and
- (ii) there are only a finite number of points at which  $F$  has a jump bigger than  $\varepsilon/2$ . From these facts it follows that for every  $\varepsilon > 0$ , there exists a constant  $C$  and another constant  $\alpha$  such that

$$P\left[\sup_x |F_n(x) - F(x)| > \varepsilon\right] \leq C \exp(-\alpha n \varepsilon^2) \quad \text{for all } n.$$

[By a result due to Dvoretzky et al. [39], this probability inequality holds with  $\alpha = 2$ .] The theorem now follows from the Borel-Cantelli Lemma. The details of the proof are left as an exercise.  $\square$

### 8.3.1 Test Statistics

There are some well-known tests for the nonparametric hypothesis  $H_0: F = F_0$  (where  $F_0$  is specified continuous cdf) against  $H_1: F \neq F_0$  which are based on the random function  $\{\sqrt{n}[F_n(x) - F_0(x)]: x \in \mathbb{R}\}$ . Test statistics for two such tests are

Kolmogorov-Smirnov statistic:  $D_n = \sup_x \sqrt{n}|F_n(x) - F_0(x)|$ , and  
Cramér-von Mises statistic:  $W_n^2 = n \int_{-\infty}^{\infty} [F_n(x) - F_0(x)]^2 dF_0(x)$ .

An important fact is the *distribution-free property* of these test statistics. To see this, suppose that  $F_0$  is *strictly increasing* (a condition which can be removed if we argue a little more carefully), and let  $U_i = F_0(X_i)$ . In [Section 1.11](#) we have seen that  $X_1, \dots, X_n$  are iid with cdf  $F_0 \iff U_1, \dots, U_n$  are iid  $\text{Unif}(0, 1)$ , so that

$$\{\sqrt{n}[F_n(x) - F_0(x)]: -\infty < x < \infty\} \stackrel{D}{=} \{\sqrt{n}[G_n \circ F_0(x) - F_0(x)], -\infty < x < \infty\}, \quad (11)$$

where  $G_n(t) = n^{-1} \sum_{i=1}^n I_{(0,t]}(U_i)$  is the edf of  $(U_1, \dots, U_n)$ .

Now letting  $t = F_0(x)$ , the statistics  $D_n$  and  $W_n^2$  can be expressed as

$$D_n = \sup_{0 \leq t \leq 1} \sqrt{n}|G_n(t) - t| \text{ and } W_n^2 = n \int_0^1 [G_n(t) - t]^2 dt. \quad (12)$$

This shows that the distribution of  $D_n$  and  $W_n^2$  under  $H_0: F = F_0$  is the same for all  $F_0$ .

The statistic  $D_n$  can also be used to construct a confidence band for an unknown continuous cdf. Since  $D_n$  is distribution-free, we can find a constant  $c_\alpha$  for a given  $0 < \alpha < 1$  such that for all  $F$  we have

$$\begin{aligned} 1 - \alpha &= P_F \left[ \sup_x \sqrt{n}|F_n(x) - F_0(x)| \leq c_\alpha \right] \\ &= P_F \left[ F_n(x) - n^{-1/2}c_\alpha \leq F(x) \leq F_n(x) + n^{-1/2}c_\alpha \text{ for all } x \right]. \end{aligned}$$

Consequently,  $\{F_n(x) \pm n^{-1/2}c_\alpha, -\infty < x < \infty\}$  is a confidence band for  $F$  with confidence coefficient  $1 - \alpha$ .

One-sided versions of the test statistic  $D_n$ , namely,

$$D_n^+ = \sup_x \sqrt{n}[F_n(x) - F_0(x)] \text{ and } D_n^- = \sup_x \sqrt{n}[F_0(x) - F_n(x)] \quad (13)$$

can be used to test  $H_0: F = F_0$  against one-sided alternatives.

The above ideas extend to two-sample problems as follows.

Let  $X_{11}, \dots, X_{1m}, \dots$  and  $X_{21}, \dots, X_{2n}, \dots$  be two independent iid sequences with common cdf's  $F_1, F_2$ , respectively, and let

$$F_{1m}(x) = m^{-1} \sum_{i=1}^m I_{(-\infty, x]}(X_{1i}) \text{ and } F_{2n}(x) = n^{-1} \sum_{i=1}^n I_{(-\infty, x]}(X_{2i})$$

be the edf based on  $X_{11}, \dots, X_{1m}$  and  $X_{21}, \dots, X_{2n}$ , respectively. Then the two-sample extensions of the Kolmogorov-Smirnov and the Cramér-von Mises statistics are, respectively,

$$\begin{aligned} D_{mn} &= \sqrt{\frac{mn}{m+n}} \sup_x |F_{1m}(x) - F_{2n}(x)| \text{ and} \\ W_{mn}^2 &= \frac{mn}{m+n} \int_{-\infty}^{\infty} [F_{1m}(x) - F_{2n}(x)]^2 d\left(\frac{mF_{1m}(x) + nF_{2n}(x)}{m+n}\right), \end{aligned} \quad (14)$$

which can be used to test  $H_0: F_1 = F_2$  vs  $H_1: F_1 \neq F_2$ .

We now examine the distribution of the stochastic process  $\{\sqrt{n}[F_n(x) - F_0(x)], -\infty < x < \infty\}$  under  $H_0: F = F_0$  or equivalently, the process  $\{\sqrt{n}[G_n(t) - t], 0 \leq t \leq 1\}$  where  $G_n(t)$  is the edf of a random sample  $U_1, \dots, U_n$  from  $\text{Unif}(0, 1)$  given in Eq. (11). Although exact results for small sample are available for the statistics  $D_n^+$  and  $D_n$ , we shall only give an outline of the asymptotic property of  $\{\sqrt{n}[G_n(t) - t], 0 \leq t \leq 1\}$  as  $n \rightarrow \infty$ .

First note that

$$\begin{aligned}\sqrt{n}[G_n(t) - t] &= n^{-1/2} \sum_{i=1}^n [I_{[0,t]}(U_i) - t], \text{ and} \\ \mathbb{E}[I_{[0,t]}(U_i)] &= t, \quad \text{Var}[I_{[0,t]}(U_i)] = t(1-t), \\ \text{Cov}[I_{[0,s]}(U_i), I_{[0,t]}(U_i)] &= \min(s, t) - st = s(1-t) \quad \text{for } 0 \leq s \leq t \leq 1.\end{aligned}$$

By multivariate CLT, it follows that for all sets of finite points  $0 < t_1 < \dots < t_k < 1$ ,

$$\begin{aligned}(\sqrt{n}[G_n(t_1) - t_1], \dots, \sqrt{n}[G_n(t_k) - t_k])^T &\xrightarrow{\mathcal{L}} N_k(\mathbf{0}, \Sigma_{t_1, \dots, t_k}), \text{ where} \\ \Sigma_{t_1, \dots, t_k} &= \begin{bmatrix} t_1(1-t_1) & t_1(1-t_2) & \cdots & t_1(1-t_k) \\ t_1(1-t_2) & t_2(1-t_2) & \cdots & t_2(1-t_k) \\ \vdots & \vdots & \vdots & \vdots \\ t_1(1-t_k) & t_2(1-t_k) & \cdots & t_k(1-t_k) \end{bmatrix}. \end{aligned} \tag{15}$$

From this, it seems plausible that the entire stochastic process  $\{\sqrt{n}[G_n(t) - t], 0 \leq t \leq 1\}$  should converge (in some sense) to a *Gaussian process* with mean value function 0 and covariance function  $\rho(s, t) = s(1-t)$ ,  $0 \leq s \leq t \leq 1$ , that is, a stochastic process  $\{X(t): 0 \leq t \leq 1\}$  such that for any  $0 \leq t_1 < \dots < t_k \leq 1$ ,  $(X(t_1), \dots, X(t_k))$  follows a  $k$ -dim normal distribution with mean vector  $\mathbf{0}$  and covariance matrix  $\Sigma_{t_1, \dots, t_k}$  given by Eq. (15). This was conjectured by Doob [40]. For a formal description of this phenomenon, we shall first review some basic facts about Brownian Motion and Weak Convergence.

### 8.3.2 Brownian Motion: Some Basic Facts

A standard Brownian Motion (B.M.) is a stochastic process  $\{X(t), t \geq 0\}$  (ie, a collection of rv's on some probability space) with the following properties:

- (i)  $X(0) = 0$  w.p. 1.
- (ii) For all  $k$  and  $0 = t_0 < t_1 < \dots < t_k$ ,  $X(t_i) - X(t_{i-1})$ ,  $1 \leq i \leq k$ , are mutually independent and  $X(t_i) - X(t_{i-1}) \sim N(0, t_i - t_{i-1})$ . Equivalently,  $(X(t_1), \dots, X(t_k))^T \sim N_n(0, ((\min(t_i, t_j))))$ . [This is the property of independent and stationary increments in a Gaussian Process.]

Wiener showed that it is possible to construct a probability distribution on a suitable  $\sigma$ -field of continuous functions on  $\mathbb{R}$  such that  $\{X(t) = X(t, \omega), t \geq 0\}$  defined on continuous functions  $\omega$  would have properties (i) and (ii). For this reason, the probability distribution of a standard B.M. is called the “Wiener Measure.” From now on, we assume that the sample paths of a B.M. are continuous w.p. 1.

In general a B.M.  $\{Y(t), t \geq 0\}$  with mean  $\mu$  and variance  $\sigma^2$  per unit time is  $Y(t) = \mu t + \sigma X(t)$ , where  $\{X(t), t \geq 0\}$  is a standard B.M.

We now state some properties of a standard B.M.

- (I) If  $\{X(t), t \geq 0\}$  is a standard B.M., then so are
  - (a)  $\{-X(t), t \geq 0\}$  (Symmetry).
  - (b)  $\{X(s+t) - X(s), t \geq 0\}$  for fixed  $s$ . Moreover, the process is independent of  $\{X(\tau), 0 \leq \tau \leq s\}$  (Markov Property).
  - (c)  $\{tX(1/t), t \geq 0\}$  (Inversion).
  - (d)  $\{\alpha^{-1/2}X(\alpha t), t \geq 0\}$  for  $\alpha > 0$  (Scale Change).
  - (e)  $\{X(t_0) - X(t_0 - t), 0 \leq t \leq t_0\}$  (Time Reversal).

For proofs of the above, it is enough to check that the covariance function of the processes in (a)–(e) is  $\min(t_1, t_2)$ .

The next three properties involve the concept of martingales. In Section 3.3, we have mentioned the martingale property of the stochastic process

$\{S_k = \sum_{i=1}^k X_i, k = 1, 2, \dots\}$  where  $\{X_i\}$  is a sequence of independent rv's with mean zero. More generally, a stochastic process  $\{X(t), t \geq 0\}$  on a probability space  $(\Omega, \mathcal{A}, P)$  is a martingale if  $E[|X(t)|] < \infty$  for all  $t$  and  $E[X(t)|\{X(\tau), 0 \leq \tau \leq s\}] = X(s)$  for all  $s \leq t$ .

- (II) (a) From the Markov Property, it follows that for  $s < t$ ,

$$\begin{aligned} E[X(t)|X(\tau), \tau \leq s] &= E[X(t) - X(s) + X(s)|X(\tau), \tau \leq s] \\ &= E[X(t) - X(s)] + X(s) = X(s), \end{aligned}$$

that is,  $\{X(t), t \geq 0\}$  is a martingale.

(b) It is also easy to verify that  $\{X(t)^2 - t, t \geq 0\}$  is a martingale.

(c) Let  $\xi(t) = e^{\theta X(t)} / E[e^{\theta X(t)}] = e^{\theta X(t) - \theta^2 t/2}$ , since

$$E[e^{\theta X(t)}] = \text{mgf of } N(0, t) \text{ at } \theta = e^{\theta^2 t/2}.$$

For  $s \leq t$

$$\begin{aligned} E[\xi(t)|\xi(\tau), \tau \leq s] &= E[\exp(\theta X(s) + \theta(X(t) - X(s)) - \theta^2 t/2)|\xi(\tau), \tau \leq s] \\ &= \exp(\theta X(s) - \theta^2 t/2) E[\exp(\theta(X(t) - X(s)))|\xi(\tau), \tau \leq s] \\ &= \exp(\theta X(s) - \theta^2 t/2) \exp(\theta^2(t-s)/2) \\ &= \exp(\theta X(s) - \theta^2 s/2) = \xi(s), \end{aligned}$$

showing that  $\{\xi(t), t \geq 0\}$  is a martingale.

- (III) Stopping Time.

For a continuous time stochastic process  $\{X(t), t \geq 0\}$ , let  $T \geq 0$  be an rv defined in such a manner that the event  $\{T \leq t\}$  depends only on  $\{X(\tau), \tau \leq t\}$ . Such an rv is called a stopping time. Examples are  $T_a = \inf\{t: X(t) = a\}$ ,  $T_{ab} = \inf\{t: X(t) \notin (a, b)\}$  for  $a < 0 < b$ .

If  $T$  is a stopping time, then  $\min(T, t)$  is also a stopping time. Suppose that  $\{X(t), t \geq 0\}$  is a martingale with continuous sample paths and let  $T < \infty$  be a stopping time with  $E[|X(T)|] < \infty$ . Then under integrability conditions,  $E[X(T)] = E[X(0)]$ . See Freedman [41] Brownian Motion and Diffusion, p. 193 and Breiman [42] Probability, Generalization of Corollary 5.31, p. 98 and 274 for details.

The martingale  $\{\xi(t), t \geq 0\}$  defined in II(c) above satisfies these conditions, so

$$E[\xi(\min(T, t))] = E[\xi(0)]. \quad (16)$$

### 8.3.3 Weak Convergence of $\{Y_n(t) = \sqrt{n}(G_n(t) - t), 0 \leq t \leq 1\}$

The random functions  $Y_n(\cdot)$  are not in  $C[0, 1]$  due to jumps in  $G_n(t)$  given in Eq. (11). However, we can take care of this difficulty by a minor adjustment of  $G_n(t)$ . Let  $U_{n:1} < \dots < U_{n:n}$  denote the order statistics of the random sample  $(U_1, \dots, U_n)$  from  $Unif(0, 1)$ . The edf  $G_n$  has jumps of  $1/n$  at each of the order statistics. Now define  $\tilde{G}_n$  by linear interpolation between the points:

$$(0, 0), (U_{n:1}, 1/(n+1)), (U_{n:2}, 2/(n+1)), \dots, (U_{n:n}, n/(n+1)), (1, 1), \text{ ie,}$$

$$\tilde{G}_n(t) = (n+1)^{-1}[(i-1) + (t - U_{n:i-1})/(U_{n:i} - U_{n:i-1})] \text{ for } U_{n:i-1} \leq t \leq U_{n:i}$$

letting  $U_{n:0} = 0$  and  $U_{n:n+1} = 1$ . Then  $|\tilde{G}_n(t) - G_n(t)| \leq 1/n$  for all  $t \in [0, 1]$ , and  $\{\tilde{Y}_n(t) = \sqrt{n}(\tilde{G}_n(t) - t), 0 \leq t \leq 1\}$  is in  $C[0, 1]$ . We can, therefore, use the theory of weak convergence in  $C[0, 1]$  outlined in Section A.5 to find the weak limit of  $\{Y_n(\cdot)\}$  which is the same as the weak limit of  $\{\tilde{Y}_n(\cdot)\}$ , because  $\sup_{0 \leq t \leq 1} |Y_n(t) - \tilde{Y}_n(t)| \leq n^{-1/2}$ .

As mentioned at the end of Section 8.3.1, the natural candidate for the weak limit of  $\{Y_n(t) = \sqrt{n}(G_n(t) - t), 0 \leq t \leq 1\}$  is a Gaussian process on  $[0, 1]$  with mean value function 0 and covariance function  $\{\rho(s, t) = s(1-t), 0 \leq s \leq t \leq 1\}$ . Now the process  $\{Y(t) = X(t) - tX(1), 0 \leq t \leq 1\}$  where  $\{X(t)\}$  is a standard B.M. fits this description, because  $\{Y(t)\}$  is a Gaussian process with  $E[Y(t)] = 0$  for all  $t$  and for  $s < t$ ,

$$\begin{aligned} \text{Cov}[Y(s), Y(t)] &= \text{Cov}[X(s), X(t)] - s\text{Cov}[X(t), X(1)] - t\text{Cov}[X(s), X(1)] + st\text{Var}[X(1)] \\ &= s - st - st + st = s(1-t). \end{aligned}$$

To accomplish the actual proof of  $Y_n(\cdot) \xrightarrow{w} Y(\cdot)$  via weak convergence of  $\{\tilde{Y}_n(\cdot)\}$ , we invoke Theorem A.5.2 in Section A.5, of which Condition (i) regarding convergence of fdd's is already seen, Condition (ii) that  $Y_n(0) = O_P(1)$  is trivial and Condition (iib) is verified by a lengthy analysis (see [43, p. 105–108]).

The process  $\{Y(t) = X(t) - tX(1), 0 \leq t \leq 1\}$  where  $\{X(t)\}$  is a standard B.M. is called a *Brownian Bridge* because it connects the points  $(0, Y(0)) = (0, 0)$  and  $(1, Y(1)) = (1, 0)$  by a continuous sample path.

### 8.3.4 Asymptotic Distributions of $D_n^+$ , $D_n$ , and $W_n^2$

All three statistics

$$\begin{aligned} D_n^+ &= \sup_{0 \leq t \leq 1} \sqrt{n}[G_n(t) - t] = \sup_{0 \leq t \leq 1} Y_n(t) \\ D_n &= \sup_{0 \leq t \leq 1} |Y_n(t)| \text{ and } W_n^2 = \int_0^1 Y_n(t)^2 dt \end{aligned}$$

given by Eqs. (12) and (13) are continuous functions of  $\{Y_n(\cdot)\}$ . It therefore follows from the Continuous Mapping [Theorem A.5.1, Section A.5](#) that

**Theorem 8.3.2.**  $D_n^+ \xrightarrow{\mathcal{L}} \sup_{0 \leq t \leq 1} Y(t)$ ,  $D_n \xrightarrow{\mathcal{L}} \sup_{0 \leq t \leq 1} |Y(t)|$ , and  $W_n^2 \xrightarrow{\mathcal{L}} \int_0^1 Y(t)^2 dt$ , where  $\{Y(t)\}$  is the Brownian Bridge on  $[0, 1]$ .

We first find the asymptotic distribution of  $D_n^+$ .

**Theorem 8.3.3.** Let  $\{Y(t) = X(t) - tX(1), 0 \leq t \leq 1\}$ , where  $\{X(t)\}$  is a standard B.M. Then

$$P\left[\sup_{0 \leq t \leq 1} Y(t) \geq y\right] = e^{-2y^2}, \quad y > 0.$$

*Proof.* Note that

$$\begin{aligned} \sup_{0 \leq t \leq 1} Y(t) \geq y &\iff \sup_{s \geq 0} Y(s/(1+s)) \geq y \\ &\iff X^*(s) = (1+s)Y(s/(1+s)) \geq (1+s)y \end{aligned}$$

for some  $s \geq 0$ . Now  $\{X^*(t) = (1+t)Y(t/(1+t)), t \geq 0\}$  is a standard B.M., because it is a Gaussian process with  $E[X^*(t)] = 0$  and

$$\text{Cov}[X^*(s), X^*(t)] = (1+s)(1+t)[s/(1+s)][1-t/(1-t)] = s, \quad \text{for } s \leq t.$$

Hence

$$P\left[\sup_{0 \leq t \leq 1} Y(t) \geq y\right] = P[X(s) \geq ys + y \quad \text{for some } s \geq 0], \tag{17}$$

where  $\{X(t)\}$  is a standard B.M.

We shall now find  $P[X(t) \geq at + b \text{ for some } t \geq 0]$  by a martingale approach.

Consider the stopping time

$$\begin{aligned} T &= \min\{t: X(t) = at + b\} \text{ if } X(t) = at + b \text{ for some } t \geq 0 \text{ and} \\ T &= \infty \text{ if } X(t) < at + b \text{ for all } t > 0. \end{aligned}$$

Then

$$P[X(t) \geq at + b \text{ for some } t \geq 0] = P[T < \infty]. \tag{18}$$

Recall that for each  $t$ ,  $T \wedge t = \min(T, t)$  is a stopping time, and  $\{\xi(t) = e^{\theta X(t)} / E[e^{\theta X(t)}] = e^{\theta X(t) - \theta^2 t/2}\}$  is a martingale (property II(c) of a standard B.M.). Consequently, Eq. (16) holds.

For  $\theta = 2a > 0$ ,

$$\begin{aligned}\xi(T \wedge t) &= \exp\left[\theta X(T \wedge t) - \theta^2(T \wedge t)/2\right] \\ &= \begin{cases} \exp\left[\theta X(t) - \theta^2 t/2\right] < \exp\left[\theta(aT + b) - \theta^2 t/2\right], & t < T \\ \exp\left[\theta X(T) - \theta^2 T/2\right] = \exp\left[\theta(aT + b) - \theta^2 T/2\right], & t \geq T \end{cases} \\ &\leq \exp\left[\theta\{a(T \wedge t) + b\} - \theta^2(T \wedge t)/2\right] \\ &= \exp\left[2a\{a(T \wedge t) + b\} - (2a)^2(T \wedge t)/2\right] = e^{2ab}.\end{aligned}\quad (19)$$

Next note that

$$\xi(t) = \exp\left[\theta X(t) - \theta^2 t/2\right] \rightarrow 0 \text{ a.s. as } t \rightarrow \infty, \quad (20)$$

because if  $\{Z_i\}$  are iid  $N(0, 1)$  and  $X(n) = \sum_{i=1}^n Z_i$ , then  $\theta X(n) - \theta^2 n/2 = \theta n(\bar{Z}_n - \theta/2) \rightarrow -\infty$  a.s.

Now consider

$$\lim_{t \rightarrow \infty} \xi(T \wedge t) = \lim_{t \rightarrow \infty} [\xi(t)I(t \leq T) + \xi(T)I(T < t)],$$

of which the first term is 0 a.s. by Eq. (20) and the second term is  $\xi(T) \lim_{t \rightarrow \infty} I(T < t) = \xi(T)I(T < \infty)$ . Thus

$$\lim_{t \rightarrow \infty} \xi(T \wedge t) = \xi(T)I(T < \infty), \text{ a.s.}$$

Hence by dominated convergence, using Eq. (19),

$$\begin{aligned}\lim_{t \rightarrow \infty} E[\xi(T \wedge t)] &= E[\xi(T)I(T < \infty)] \\ &= E\left[\exp[\theta(aT + b) - \theta^2 T/2]I(T < \infty)\right] \\ &= E\left[\exp[2a(aT + b) - (2a)^2 T/2]I(T < \infty)\right] \\ &= e^{2ab}P[T < \infty]\end{aligned}$$

for  $\theta = 2a$ . But by Eq. (16),  $E[\xi(T \wedge t)] = 1$  for all  $t$ , so

$$\lim_{t \rightarrow \infty} E[\xi(T \wedge t)] = 1 = e^{2ab}P[T < \infty].$$

Thus  $P[T < \infty] = e^{-2ab}$  for  $a > 0$  and  $b \geq 0$ . Taking  $a = b = y > 0$  in this formula and using Eqs. (17) and (18), the theorem follows.  $\square$

The distribution of  $\sup_{0 \leq t \leq 1} |Y(t)|$  is obtained by lengthy analysis of the joint behavior of the maximum and the minimum of a standard B.M. We state the result in the following theorem, referring to Billingsley [43, p. 77–80 and 83–85] for a proof.

**Theorem 8.3.4.** *Let  $\{Y(t), 0 \leq t \leq 1\}$  be a Brownian Bridge on  $[0, 1]$ . Then*

$$P\left[\sup_{0 \leq t} |Y(t)| \leq y\right] = 1 + 2 \sum_{k=1}^{\infty} (-1)^k e^{-2k^2 y^2}, \quad y > 0.$$

The distribution of  $\int_0^1 Y(t)^2 dt$  is obtained by using an orthogonal expansion of a standard B.M. in terms of an infinite sequence of iid  $N(0, 1)$  rv's  $\{Y_0, Y_1, \dots\}$  due to Wiener [44]. We sketch the derivation of this distribution, referring to Breiman [42], Section 12.7, and Hájek and Šidák [4], V3.3 Theorem C for details.

Let  $\{X(t)\}$  be a standard B.M. on  $[0, \pi]$ . Then  $\text{Cov}[X(s), X(t)]$  can be expressed by the following identity:

$$st/\pi + (2/\pi) \sum_{k=1}^{\infty} (\sin ks)(\sin kt)/k^2 = \min(s, t) \quad \text{for } 0 \leq s, t \leq \pi. \quad (21)$$

Hence  $\{X(t)\}$  has the representation:

$$X(t) = (1/\sqrt{\pi})tY_0 + \sqrt{2/\pi} \sum_{k=1}^{\infty} \{(\sin kt)/k\}Y_k, \quad 0 \leq t \leq \pi, \quad (22)$$

where  $Y_0, Y_1, \dots$  are iid  $N(0, 1)$  rv's, because the right-hand side of Eq. (22) is a Gaussian process with  $E[X(t)] = 0$  and  $\text{Cov}[X(s), X(t)]$  is the same as in Eq. (21).

Rewrite Eq. (21), replacing  $s/\pi, t/\pi$  by  $s^*, t^*$ , respectively, and rewrite Eq. (22) by scale change to obtain

$$X^*(t) = (1/\sqrt{\pi})X(\pi t) = tY_0 + \sqrt{2} \sum_{k=1}^{\infty} \{(\sin k\pi t)/(k\pi)\}Y_k, \quad 0 \leq t \leq 1$$

and

$$s^*t^* + 2 \sum_{k=1}^{\infty} (\sin k\pi s^*)(\sin k\pi t^*)/(k\pi)^2 = \min(s^*, t^*), \quad 0 \leq s^*, t^* \leq 1. \quad (23)$$

Hence  $\{X^*(t)\}$  is a standard B.M. on  $[0, 1]$ , because  $\text{Cov}[X^*(s), X^*(t)] = \min(s, t)$  by Eq. (23).

Since  $X^*(1) = Y_0$ , it follows that

$$Y(t) = X^*(t) - tX^*(1) = \sqrt{2} \sum_{k=1}^{\infty} \{(\sin k\pi t)/(k\pi)\}Y_k \quad (24)$$

is a Brownian Bridge on  $[0, 1]$ . We thus arrive at the following theorem.

**Theorem 8.3.5.**  $\int_0^1 Y(t)^2 dt = \sum_{k=1}^{\infty} (k\pi)^{-2} Y_k^2$ , which is a mixture of  $\chi^2$ 's.

*Proof.* By Eq. (24),

$$\int_0^1 Y(t)^2 dt = 2 \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} (jk\pi^2)^{-1} Y_j Y_k \int_0^1 \sin j\pi t \sin k\pi t dt = \sum_{k=1}^{\infty} (k\pi)^{-2} Y_k^2,$$

since  $2 \int_0^1 (\sin j\pi t)(\sin k\pi t) dt = 1$  when  $j = k$  and = 0 when  $j \neq k$ .

□

Finally, we consider the two-sample statistics  $D_{mn}$  and  $W_{mn}^2$  given by Eq. (14) under the null hypothesis  $F_1 = F_2 = F$  (unknown).

For  $m, n \rightarrow \infty$  in such a way that  $m/(m+n) = m/N \rightarrow \lambda$  and  $n/(m+n) = n/N \rightarrow 1-\lambda$  for  $0 < \lambda < 1$ , write

$$\begin{aligned} D_{mn} &= \sqrt{\frac{mn}{m+n}} \sup_x |\{F_{1m}(x) - F(x)\} - \{F_{2n}(x) - F(x)\}| \\ &= \sqrt{\frac{mn}{m+n}} \sup_{0 \leq t \leq 1} |m^{-1/2} [\sqrt{m}\{G_{1m}(t) - t\}] - n^{-1/2} [\sqrt{n}\{G_{2n}(t) - t\}]| \\ &\xrightarrow{\mathcal{L}} \sqrt{\lambda(1-\lambda)} \sup_{0 \leq t \leq 1} |\lambda^{-1/2} Y_1(t) - (1-\lambda)^{-1/2} Y_2(t)| \\ &= \sup_{0 \leq t \leq 1} |\sqrt{1-\lambda} Y_1(t) - \sqrt{\lambda} Y_2(t)|, \end{aligned}$$

where  $\{Y_1(t)\}$  and  $\{Y_2(t)\}$  are independent Brownian Bridges on  $[0, 1]$ . Since  $\{Y(t) = \sqrt{1-\lambda} Y_1(t) - \sqrt{\lambda} Y_2(t)\}$  is a Gaussian process with mean value function 0 and

$$\text{Cov}[Y(s), Y(t)] = (1-\lambda)s(1-t) + \lambda s(1-t) = s(1-t) \quad \text{for } 0 \leq s, t \leq 1,$$

$D_{mn} \xrightarrow{\mathcal{L}} \sup_{0 \leq t \leq 1} |Y(t)|$  and  $\{Y(t)\}$  is a Brownian Bridge on  $[0, 1]$ . Thus the asymptotic distribution of  $D_{mn}$  is the same as that of the one-sample statistic  $D_n$  under null hypothesis.

For  $W_{mn}^2$ , analogous argument holds with the additional observation that  $(m+n)^{-1} [mG_{1m}(t) + nG_{2n}(t)] \rightarrow t$ , a.s. uniformly by Glivenko-Cantelli Lemma. Hence  $W_{mn}^2 \xrightarrow{\mathcal{L}} \int_0^1 Y(t)^2 dt$ .

## Exercises

- 8.1.** Let  $p$  be the joint pdf of  $X_1, \dots, X_m, X_{m+1}, \dots, X_{m+n}$ . We want a UMP similar region test at level  $\alpha$  for  $H_0: p \in P_0$  vs  $H_1: p = q_\Delta$ ,  $\Delta > 1$ , where

$$q_\Delta(x_1, \dots, x_{n+m}) = (2\pi)^{(m+n)/2} \Delta^{-n/2} \exp \left[ -(1/2) \sum_{i=1}^m x_i^2 - (2\Delta)^{-1} \sum_{i=m+1}^{n+m} x_i^2 \right].$$

- (a)** Show that the UMP similar region test at level  $\alpha$  for  $H_0$  vs  $H_1$  is of the form:  
 $\varphi(\mathbf{x}_{(m+n)}, \mathbf{r}_{(m+n)}) = 0$ , or  $\gamma(\mathbf{x}_{(m+n)})$ , or 1  
according as  $\sum_{i=m+1}^{m+n} x_i^2 <$ , or  $=$ , or  $> k(\mathbf{x}_{(m+n)})$ .
- (b)** For  $m = 6$ ,  $n = 9$ , and  $\Delta = 2$ , generate  $(x_1, \dots, x_{m+n})$  from  $q_\Delta$  and apply the above test at level  $\alpha = 0.05$  by determining  $k(x_{m+n})$  and  $\gamma(x_{m+n})$  from  
(i) the permutation distribution of  $\sum_{i=m+1}^{m+n} x_i^2$ ,  
(ii) a random sample of 100 combinations.

**8.2.** The Mann-Whitney statistic, properly normalized, is

$$T_n = \frac{\sqrt{n}}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \{I_{[0,\infty)}(X_i - Y_j) - \theta\},$$

$$n_1 = [n\lambda], n_2 = [n(1-\lambda)], \quad 0 < \lambda < 1,$$

where  $\{X_i\}$  and  $\{Y_j\}$  are independent iid sequences from  $F$  and  $G$ , respectively, and  $\theta = \int G(x) dF(x)$ . Find

- (a) the Hájek projection  $T_n^*$  of  $T_n$  when  $G = F$ ,
  - (b) the asymptotic distribution of  $T_n^*$ , and
  - (c) the asymptotic distribution of  $T_n$  with proper justification.
- 8.3.** The Kendall's tau statistic based on a random sample  $(X_1, Y_1), \dots, (X_N, Y_N)$  from a continuous bivariate distribution is

$$\tau_n = \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j \neq i}^N sign(X_i - X_j) sign(Y_i - Y_j).$$

- (a) Find the Hájek projection  $\tau_N^*$  of  $\tau_N$  and the asymptotic distribution of  $\sqrt{N}\tau_N^*$  under  $H_0: X$  and  $Y$  are independent.
  - (b) Use (a) and the property of  $\text{Var}[\tau_N - \tau_N^*]$  to derive the asymptotic distribution of  $\sqrt{N}\tau_N$  under  $H_0$ .
- 8.4.** Let  $X_1, \dots, X_m, Y_1 = X_{m+1}, \dots, Y_n = X_{m+n}$  be independent rv's,  $N = m + n$ .
- (a) Under  $H_0: X_i \sim N(0, 1)$ ,  $i = 1, \dots, N$  and under contiguous alternative  $H_N: X_i \sim N(0, 1)$ ,  $i = 1, \dots, m$  and  $X_i \sim N(\delta/\sqrt{N}, 1)$ ,  $i = m+1, \dots, N$ , find the asymptotic distributions of the following:

$$S_{N1} = \text{Two-sample } t\text{-statistic} = (\bar{X} - \bar{Y}) \left/ \left[ \left( \frac{1}{m} + \frac{1}{n} \right) \frac{(m-1)S_X^2 + (n-1)S_Y^2}{m+n-2} \right]^{1/2} \right.,$$

$$S_{N2} = \text{Wilcoxon statistic} = m^{-1} \sum_{i=1}^m R_{N:i} - n^{-1} \sum_{i=m+1}^{m+n} R_{N:i}, \text{ or equivalently,}$$

$$S_{N2} = \sum_{i=1}^m R_{N:i},$$

$$S_{N3} = \text{Fisher-Yates normal scores rank statistic} = \sum_{i=1}^m E[\Phi^{-1}(U_{N:R_{N:i}})], \text{ using usual notations.}$$

- (b) Find the Pitman AREs of  $S_{N2}$  and  $S_{N3}$  with respect to  $S_{N1}$ .

- 8.5.** Let  $X_1, \dots, X_n$  be a random sample from a distribution with mean  $\mu$ , variance  $\sigma^2$ , and finite fourth central moment  $\beta = E[(X - \mu)^4]$ . We want the asymptotic distribution of the sample variance  $s_n^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ .
- (a) Show that  $s_n^2 = \{n(n-1)\}^{-1} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)^2$ , which is a  $U$ -statistic.
  - (b) Use this  $U$ -statistic form to find the Hájek projection of  $2s_n^2$  and then the asymptotic distribution of  $\sqrt{n}(s_n^2 - \sigma^2)$ .
- 8.6.** Let  $(X_1, Y_1), \dots, (X_N, Y_N)$  be iid as  $(X, Y)$  following the bivariate normal distribution with  $E[X] = E[Y] = 0$ ,  $E[X^2] = E[Y^2] = 1$ , and  $E[XY] = \theta$ . For testing  $H_0: \theta = 0$  vs  $H_1: \theta > 0$  (ie, independence vs positive dependence), find the Pitman ARE of

$S_{N1} = \sqrt{N}\tau_N$  with respect to  $S_{N2} = \sqrt{N}r_N$ , where  $\tau_N$  is Kendall's tau and  $r_N$  is the product-moment correlation based on  $(X_1, Y_1), \dots, (X_N, Y_N)$  as in [Example 8.1.3](#). Use the Hájek projection of  $\tau_N$  and its asymptotic property obtained in Exercise 8.3 and show that  $\sqrt{N}r_N = N^{-1/2} \sum_{i=1}^N X_i Y_i + o_P(1)$ . Now find  $\dot{l}(x, y; 0)$  and  $I$ , express  $\log L_N$ ,  $S_{N1}$ , and  $S_{N2}$  in the desired form. [Hint: If  $\Phi$  and  $\phi$  are the cdf and pdf of  $N(0, 1)$ , then  $\int_{-\infty}^{\infty} x\Phi(x) d\Phi(x) = \int_{-\infty}^{\infty} [\int_{-\infty}^x \phi(t) dt] x\phi(x) dx = \int_{-\infty}^{\infty} [\int_t^{\infty} x\phi(x) dx] \phi(t) dt$ . Evaluate this.]

- 8.7.** In this problem, Kendall's tau and Spearman's rank correlation are used in a different context. Let  $X_i = \theta d_i + Z_i$ ,  $i = 1, \dots, n$ , where  $Z_1, \dots, Z_n$  are iid with a continuous cdf and  $d_1, \dots, d_n$  are equally spaced constants in increasing order, which can be taken to be  $1, \dots, n$ , without loss of generality. The following statistics can be used to test  $H_0: \theta = 0$  vs  $H_1: \theta > 0$ .

(i) Moore-Wallis Difference-Sign Statistic:  $D_n = \sum_{i=2}^n I_{(0,\infty)}(X_i - X_{i-1})$ ,

(ii) Difference-Sign Correlation Coefficient (Kendall):

$$\tau_N = [4n(n-1)]^{-1} \sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{(0,\infty)}(X_i - X_j),$$

(iii) Rank Correlation Coefficient (Spearman):

$$\rho_n = 1 - [6n(n^2 - 1)]^{-1} \sum_{i=1}^{n-1} \sum_{j=i+1}^n (j - i)[1 + sign(X_i - X_j)],$$

(iv) The  $t$ -statistic:  $T_n = \sum_{i=1}^n (d_i - \bar{d}_n)(X_i - \bar{X}_n) / \sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2}$ .

We now ask the following questions, regarding the asymptotic properties of these statistics:

- (a) Express Kendall's tau and Spearman's rank correlation coefficient in the form given in (ii) and (iii).
- (b) Express the  $t$ -statistic in the form given in (iv). [This is the test statistic for  $H_0: \theta = 0$  in the linear model  $X_i = \mu + d_i\theta + Z_i$ , where  $Z_1, \dots, Z_n$  are iid  $N(0, \sigma^2)$ , obtained from the least squares estimator and the residual sum of squares. See [Example 6.9.10](#).]
- (c) Find the asymptotic distribution of  $D_n$ ,  $\tau_n$ ,  $\rho_n$ , and  $T_n$  (properly normalized) under  $H_0$  and under  $H_1$ , assuming that the  $Z_i$ 's are iid as  $N(0, \sigma^2)$ . For  $H_1$ , consider contiguous alternatives.
- (d) Calculate the Pitman AREs of the tests based on  $D_n$ ,  $\tau_n$ , and  $\rho_n$  with respect to the test based on  $T_n$ .

[Hint: The statistic  $\tau_n$  is a  $U$ -statistic, but  $\rho_n$  is not; so find the Hájek projection (adjusting for the mean) and check that the Hájek projection differs from the

original statistic by  $o_P(1)$ . Express  $T_n$  in terms of  $b_n = \sqrt{\sum_{i=1}^n (d_i - \bar{d}_n)^2}$ ,

$Y_n = \sum_{i=1}^n (d_i - \bar{d}_n) Z_i / b_n \sim N(0, 1)$ , and  $W_n = \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 - Y_n^2 \sim \chi_{n-2}^2$ , and in the nonnull case use  $\theta = \theta_n = \delta/b_n$ . For  $D_n$ , use [Theorem 3.3.3](#) of [Chapter 3](#) for  $m$ -dependent processes.]

- 8.8.** Let  $U_{n1} = (1/n^{(r)}) \sum_{n,r} g_1(X_{i_1}, \dots, X_{i_r})$  and  $U_{n2} = (1/n^{(s)}) \sum_{n,s} g_2(X_{i_1}, \dots, X_{i_s})$  be two  $U$ -statistics based on iid rv's  $X_1, \dots, X_n$ , both  $g_1$  and  $g_2$  being symmetric in their arguments. Let  $\theta_1 = E[g_1(X_1, \dots, X_r)]$ ,  $\theta_2 = E[g_2(X_1, \dots, X_s)]$  and assume that

$E[g_1^2(X_1, \dots, X_r)]$ , and  $E[g_2^2(X_1, \dots, X_s)]$  are finite. We want the asymptotic distribution of  $\sqrt{n}[\varphi(U_{n1}, U_{n2}) - \varphi(\theta_1, \theta_2)]$ , where  $\varphi$  has continuous first partial derivatives at  $(\theta_1, \theta_2)$ . Let

$$\begin{aligned} h_1(x_1, \dots, x_r) &= g_1(x_1, \dots, x_r) - \theta_1, \quad h_2(x_1, \dots, x_s) = g_2(x_1, \dots, x_s) - \theta_2, \\ h_1^*(X_1) &= E[h_1(X_1, \dots, X_r)|X_1], \quad h_2^*(X_1) = E[h_2(X_1, \dots, X_s)|X_1], \\ \sigma_{11} &= E[h_1^{*2}(X_1)], \quad \sigma_{22} = E[h_2^{*2}(X_1)], \quad \sigma_{12} = E[h_1^*(X_1)h_2^*(X_1)], \text{ and} \end{aligned}$$

suppose that  $\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix}$  is positive-definite.

- (a) Find two-dimensional Hájek projection  $(V_{n1}^*, V_{n2}^*)$  of  $(V_{n1}, V_{n2}) = (\sqrt{n}(U_{n1} - \theta_1), \sqrt{n}(U_{n2} - \theta_2))$ .
- (b) Show that  $(V_{n1}^*, V_{n2}^*) - (V_{n1}, V_{n2}) = o_P(1)$ .
- (c) Find the asymptotic distribution of  $(V_{n1}, V_{n2})$ .
- (d) From (c), find the asymptotic distribution of  $\sqrt{n}[\varphi(U_{n1}, U_{n2}) - \varphi(\theta_1, \theta_2)]$ , using the delta method.

- 8.9.** Let  $X_1, \dots, X_n$  be iid positive-valued rv's with cdf  $F$  having mean  $\mu$  and variance  $\sigma^2$ . Let

$$\begin{aligned} \Delta &= E_F[|X_1 - X_2|] = \iint |x_1 - x_2| dF(x_1) dF(x_2) \text{ and} \\ \zeta &= E_F[X_1|X_1 - X_2|] = \iint x_1|x_1 - x_2| dF(x_1) dF(x_2). \end{aligned}$$

The statistic  $G_n = D_n/(2\bar{X}_n)$ , where  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$  and  $D_n = [n(n-1)]^{-1} \sum_{i=1}^n \sum_{j \neq i=1}^n |X_i - X_j|$  is known as Gini's coefficient of concentration.

Use the result of Exercise 8.8 to find the asymptotic distribution of  $\sqrt{n}(G_n - \Delta/(2\mu))$ .

- 8.10.** Verify that the score functions

$$\begin{aligned} a_N^+(i, f) &= E_{q_0}\left[-\frac{f'}{f}(|X|_{N:i})\right] = E[\varphi^+(U_{N:i}, f)] \text{ and} \\ a_{1N}^+(i, f) &= E_{q_0}\left[-1 - |X|_{N:i}\frac{f'}{f}(|X|_{N:i})\right] = E[\varphi_1^+(U_{N:i}, f)] \end{aligned}$$

are as stated in [Theorems 8.2.2 and 8.2.3](#).

- 8.11.** Verify that the scores and the resulting test statistics for the specific problems in [Section 8.2.5](#) are as stated in the text.
- 8.12.** Let  $p$  be the joint pdf of  $(X_1, \dots, X_{m+n})$ . We want to test  $H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q_\Delta$ ,  $\Delta > 0$ . Show that if the joint cdf  $Q_\Delta$  corresponding to  $q_\Delta$  is given by (Lehmann Alternative)

$$Q_\Delta[X_i \leq x_i, 1 \leq i \leq m+n] = \prod_{i=1}^m \left[ (1 - \Delta)F(x_i) + \Delta F^2(x_i) \right] \prod_{i=m+1}^{m+n} F(x_i),$$

where  $F$  is an arbitrary absolutely continuous cdf, then the LMP rank test of

$H_0: p \in \mathcal{P}_0$  vs  $H_1: p = q_\Delta, \Delta > 0$ , has a critical region of the form

$$\sum_{i=1}^m R_{m+n:i} \geq \text{constant}.$$

- 8.13.** In the set-up of Exercise 8.12, suppose that the joint cdf  $Q_\Delta$  is given by

$$Q_\Delta[X_i \leq x_i, 1 \leq i \leq m+n] = \prod_{i=1}^m F(x_i)^{1+\Delta} \prod_{i=m+1}^{m+n} \left[ 1 - \{1 - F(x_i)\}^{1+\Delta} \right],$$

where  $F$  is an arbitrary absolutely continuous cdf. Show that the LMP rank test of

$H_0: \mathcal{P} \in \mathcal{P}_0$  vs  $H_1: p = q_\Delta, \Delta > 0$ , has a critical region of the form

$$\sum_{i=1}^m a_N(R_{m+n:i}) \geq \text{constant}, \text{ with } a_N(i) = \sum_{j=0}^{i-1} 1/(N-j) - \sum_{j=0}^{N-i} 1/(N-j),$$

$N = m+n$ . [Hint: First show that  $E[-\log(1 - U_{N:i})] = \sum_{j=0}^{i-1} 1/(N-j)$  and

$E[-\log U_{N:i}] = \sum_{j=0}^{N-i} 1/(N-j)$ , where  $U_{N:1} < \dots < U_{N:N}$  is an ordered random sample from  $\text{Uniform}(0, 1)$ .]

- 8.14.** Verify the formula for  $\text{Var}[B(X)]$  given in Eq. (9) and the limit of  $\lambda_N(1 - \lambda_N)^{-1} V_N$  as  $\theta_N \rightarrow 0$  given by Eq. (10).

- 8.15.** Suppose that  $X_1, \dots, X_N$  are iid rv's and let  $R_{N:i}$  be the rank of  $X_i$  among  $X_1, \dots, X_N$ . Consider the scores

$$a_N(i) = \begin{cases} i & 1 \leq i \leq N-1 \\ N^2/2 & i = N \end{cases},$$

and let  $S_{mn} = \sum_{i=1}^m a_N(R_{N:i})$ ,  $m+n=N$ . Show that if  $N \rightarrow \infty$ ,  $m \rightarrow \infty$ , and  $m/N \rightarrow 0$ , then

(a)  $\{E[S_{mn}] - mN/2\}^2 / (mnN/12) \rightarrow \infty$ ,

(b)  $\text{Var}[S_{mn}] / (mnN/12) \rightarrow \infty$ , and

(c)  $(S_{mn} - mN/2) \xrightarrow{\mathcal{L}} N(0, 1)$ .

- 8.16.** Give a detailed proof of the Glivenko-Cantelli Theorem using the outline in the text.

- 8.17.** Let  $\{X(t), 0 \leq t \leq 1\}$  be a collection of rv's such that

(i)  $X(0) = 0$  with probability 1,

(ii) for any  $0 = t_0 < t_1 < \dots < t_k \leq 1$ , the increments  $X(t_i) - X(t_{i-1})$  are independent  $N(0, t_i - t_{i-1})$ .

Show that for any  $0 < t_1 < \dots < t_k \leq 1$  ( $X(t_1), \dots, X(t_k)$ ) follows the  $k$ -dim normal distribution with mean vector  $\mathbf{0}$  and covariance matrix  $((\sigma(t_i, t_j))) = ((\min(t_i, t_j)))$ .

- 8.18.** Let  $U_1, U_2, \dots$  be iid  $\text{Uniform}(0, 1)$  rv's and  $F_n(t) = n^{-1} \sum_{i=1}^n I_{[0,t]}(U_i)$ . Let

$$Y_n(t) = \sqrt{n}[F_n(t) - t].$$

Show that for any  $0 \leq t_1 < \dots < t_k \leq 1$ ,  $(Y_n(t_1), \dots, Y_n(t_k)) \xrightarrow{\mathcal{L}} (Y(t_1), \dots, Y(t_k))$ , where  $Y(t) \stackrel{\mathcal{D}}{=} X(t) - tX(1)$  and  $\{X(t): 0 \leq t \leq 1\}$  is an in Exercise 8.17.

# Curve Estimation

## 9.1 Introduction

This chapter is concerned with three problems:

**Problem 1.** Let  $X_1, \dots, X_n$  be a random sample from a cdf  $F$  with pdf  $f = F'$ . We want to estimate the function  $f$ .

**Problem 2.** Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be a random sample from a bivariate distribution with regression function

$$\begin{aligned} m(x) &= E[Y|X = x] = \int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy = \varphi(x)/f(x), \text{ where} \\ \varphi(x) &= \int_{-\infty}^{\infty} y f_{XY}(x,y) dy \end{aligned}$$

and  $f$  is the pdf of  $X$ . We want to estimate the function  $m$ .

**Problem 3.** Let  $(T, C)$  be independent positive-valued rv's, where  $T$  is the survival time of sample unit with cdf  $F$  and pdf  $f$ , observations on which may be stopped at time  $C$ . We want to estimate the survival function  $S(t) = 1 - F(t)$  and the hazard function  $\lambda(t) = f(t)/S(t)$  based on iid observations on  $(T, C)$ .

We shall develop methods for estimating the functions  $f$  and  $m$ , and look at the asymptotic properties of these estimators as  $n \rightarrow \infty$ . Methods for estimating the survival function and integrated hazard function will also be constructed.

## 9.2 Density Estimation

Since the empirical cdf  $F_n(x) = n^{-1} \sum_{i=1}^n I_{(-\infty, x]}(X_i)$  is a natural estimator of  $F$ , we can attempt to estimate  $f = F'$  via the estimator  $F_n$  of  $F$ , using the relation

$$\begin{aligned} f(x) &= \lim_{h \downarrow 0} h^{-1} [F(x + h/2) - F(x - h/2)] \\ &\approx h^{-1} [F(x + h/2) - F(x - h/2)] \end{aligned}$$

for small  $h > 0$ . This leads to the estimator

$$\begin{aligned} f_n(x) &= h_n^{-1}[F_n(x + h_n/2) - F_n(x - h_n/2)] \\ &= (nh_n)^{-1} \sum_{i=1}^n I_{(x-h_n/2, x+h_n/2]}(X_i) \\ &= (nh_n)^{-1} \sum_{i=1}^n I_{(-1/2, 1/2]} \left( \frac{x - X_i}{h_n} \right). \end{aligned} \quad (1)$$

The data distribution is a discrete distribution with mass  $1/n$  at each  $X_i$ . If we spread the discrete mass uniformly over an interval of length  $h_n$  centered around  $X_i$ , then the mass  $1/n$  at  $X_i$  is replaced by a histogram of height  $(nh_n)^{-1}$  on the interval  $(X_i - h_n/2, X_i + h_n/2]$ . Putting all these histograms around  $X_1, \dots, X_n$  together, we obtain  $f_n(x)$  defined by Eq. (1). In this intuitive description, we could spread the mass  $1/n$  at each  $X_i$  by an arbitrary pdf  $h_n^{-1}K((\cdot - X_i)/h_n)$  instead of the uniform pdf  $h_n^{-1}I_{(-1/2, 1/2]}((\cdot - X_i)/h_n)$ . This would define a general class of estimators of the form

$$f_n(x) = (nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n), \quad (2)$$

where  $K$  is a pdf (ie,  $K(u) \geq 0$  and  $\int_{-\infty}^{\infty} K(u) du = 1$ ). Such an estimator is called a *kernel estimator* with *kernel*  $K$  and *bandwidth*  $h_n$ . Since  $h_n$  serves the purpose of spreading discrete masses over the support of  $K$ , it is called a *smoothing parameter* of the estimator  $f_n(x)$ .

Note that  $f_n$  is itself a pdf, because  $f_n(x) \geq 0$  for all  $x$  and it is easy to check that  $\int_{-\infty}^{\infty} f_n(x) dx = 1$ .

The kernel  $K$  and the bandwidth  $h_n$  are chosen by the user. The following kernels are often used:

$$\begin{aligned} \text{Uniform: } K(u) &= I_{[-1/2, 1/2]}(u), \\ \text{Logistic: } K(u) &= e^u / (1 + e^u)^2, \\ \text{Epanechnikov: } K(u) &= (3/4)(1 - u^2)I_{[-1, 1]}(u). \end{aligned}$$

Actually, the estimators are not very sensitive to the choice of  $K$ , but they are very sensitive to the choice of  $h_n$ .

For early work on kernel estimators of density functions and regression functions (discussed in the next section), we refer to Rosenblatt [45], Parzen [46], and Nadaraya [47].

### Properties of Kernel Estimators of Density Functions

We first look at the effect of  $h_n$  on the estimator  $f_n$  in terms of its mean and variance.

It should be intuitively clear that the bias increases as  $h_n$  increases (eg, with uniform kernel,  $f_n(x)$  estimates  $h_n^{-1}[F(x + h_n/2) - F(x - h_n/2)]$  instead of  $\lim_{h \downarrow 0} h^{-1}[F(x + h/2) - F(x - h/2)]$ ), and the variance increases as  $h_n$  decreases (because fewer observations

make appreciable contributions as the bandwidth shrinks). Therefore, we should choose  $h_n$  appropriately to strike a balance between bias and variance in order to achieve the best rate at which the mean-square error ( $MSE = Bias^2 + Variance$ ) tends to 0 as  $n \rightarrow \infty$ , remembering that  $h_n \downarrow 0$  as  $n \rightarrow \infty$  anyway.

$$\begin{aligned} E[f_n(x)] &= E\left[h_n^{-1}K((x-X)/h_n)\right] = \int_{-\infty}^{\infty} h_n^{-1}K((x-y)/h_n)f(y) dy \\ &= \int_{-\infty}^{\infty} K(u)f(x-h_n u) du, \end{aligned} \quad (3)$$

letting  $u = (x-y)/h_n$ , and similarly,

$$\begin{aligned} \text{Var}[f_n(x)] &= n^{-1} \left[ h_n^{-1} \int_{-\infty}^{\infty} K^2(u)f(x-h_n u) du \right. \\ &\quad \left. - \left\{ \int_{-\infty}^{\infty} K(u)f(x-h_n u) du \right\}^2 \right]. \end{aligned} \quad (4)$$

We make the following assumptions:

1. (a)  $K$  is a symmetric pdf (ie,  $K(-u) = K(u)$  for all  $u$ ), (b)  $\sigma_K^2 = \int_{-\infty}^{\infty} u^2 K(u) du < \infty$ , and  
(c)  $\|K\|^2 = \int_{-\infty}^{\infty} K^2(u) du < \infty$ .
2.  $f''$  is bounded and continuous.

Assumption 1 poses no problem, because  $K$  is chosen by the user. If Assumption 2 is replaced by other smoothness conditions of  $F$ , then the results we now derive, will change accordingly.

Expand  $f(x-h_n u)$  in the expressions (3) and (4) of  $E[f_n(x)]$  and  $\text{Var}[f_n(x)]$ :

$$\begin{aligned} f(x-h_n u) &= f(x) - h_n u f'(x) + (1/2)h_n^2 u^2 f''(x) \\ &\quad + (1/2)h_n^2 u^2 \{f''(-\lambda h_n u) - f''(x)\}, \quad 0 \leq \lambda \leq 1. \end{aligned} \quad (5)$$

Since  $\int_{-\infty}^{\infty} uK(u) du = 0$  by Assumption 1(a) and

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} \{f''(x-h_n u) - f''(x)\} u^2 K(u) du = 0$$

by Assumption 2 and Dominated Convergence, using Eq. (5) we have

$$\begin{aligned} \text{Bias}[f_n(x)] &= E[f_n(x)] - f(x) = (1/2)h_n^2 \left[ \sigma_K^2 f''(x) + o(1) \right] \text{ and} \\ \text{Var}[f_n(x)] &= (nh_n)^{-1} \left[ \|K\|^2 f(x) + o(1) \right]. \end{aligned}$$

Combining the above two formulas, we have

$$MSE[f_n(x)] = \left( h_n^4 \right) \left[ \sigma_K^4 f''(x)^2 / 4 + o(1) \right] + (nh_n)^{-1} \left[ \|K\|^2 f(x) + o(1) \right],$$

so that

$$n^{4/5} MSE[f_n(x)] = \left( n^{1/5} h_n \right)^4 \left[ \sigma_K^4 f''(x)^2 / 4 + o(1) \right] + \left( n^{1/5} h_n \right)^{-1} \left[ \|K\|^2 f(x) + o(1) \right]. \quad (6)$$

On the right-hand side of Eq. (6), the first term  $\rightarrow \infty$  if  $n^{1/5}h_n \rightarrow \infty$  and the second term  $\rightarrow \infty$  if  $n^{1/5}h_n \rightarrow 0$ . Thus  $n^{4/5}\text{MSE}[f_n(x)] \rightarrow \infty$  if  $n^{1/5}h_n$  either tends to 0 or to  $\infty$ , and remains bounded if  $h_n$  is of the order of magnitude  $n^{-1/5}$ . Taking  $h_n = tn^{-1/5}$ , we have

$$n^{4/5}\text{MSE}[f_n(x)] = t^4\sigma_K^4 f''(x)^2/4 + t^{-1}\|K\|^2 f(x) + o(1).$$

Finally,  $t^4a + t^{-1}b$  with  $a, b > 0$  is minimized at  $t_0 = (b/(4a))^{1/5}$ .

Using this, we see that  $n^{4/5}\text{MSE}[f_n(x)]$  is minimized with  $h_n = n^{-1/5}t_0$  where  $t_0 = [\{\|K\|^2 f(x)\}/\{\sigma_K^4 f''(x)^2\}]^{1/5}$ .

For example, if  $K$  is the uniform kernel on  $[-1/2, 1/2]$ , then with  $\|K\|^2 = 1$  and  $\sigma_K^2 = 1/12$ , we have  $t_0 = [144f(x)/f''(x)^2]^{1/5}$ .

Since  $t_0$  involves the unknown  $f(x)$  and  $f''(x)$ , one way to implement this in practice would be to obtain initial (consistent) estimates of  $f(x)$  and  $f''(x)$  from the data and then “plug in” these estimates in the formula for  $t_0$ .

So far we have considered the estimation of  $f(x)$  at a specific  $x$ . However, we are often interested in estimating the entire function  $f$ , in which case, one would like to use the same bandwidth for all  $x$ . For this, the integrated mean-square error  $\text{IMSE}(f_n) = \int \text{MSE}[f_n(x)] dx$  or integrated square error  $\int [f_n(x) - f(x)]^2 dx$  would be a reasonable criterion to minimize. There is a huge literature on the issue of bandwidth choice in density estimation and regression estimation, some of which will be discussed at the end of this chapter.

We shall now establish the following asymptotic properties of  $f_n$ :

**Theorem 9.2.1** (Strong Uniform Consistency). *In addition to Assumptions 1(a, b) and 2, suppose that the kernel  $K$  is of bounded variation on  $(-\infty, \infty)$ . Then  $\sup_x |f_n(x) - f(x)| \rightarrow 0$  a.s., provided that  $nh_n^2/\log n \rightarrow \infty$  as  $n \rightarrow \infty$ . [In particular, this holds for  $h_n = O(n^{-1/5})$ .]*

**Theorem 9.2.2** (Asymptotic Normality). *Under Assumptions 1(a, b, c) and 2, with  $h_n = n^{-1/5}t$ ,*

$$n^{2/5}[f_n(x) - f(x)] \xrightarrow{\mathcal{L}} N\left((1/2)t^2\sigma_K^2 f''(x), t^{-1}\|K\|^2 f(x)\right).$$

We now prove the strong uniform consistency property. The asymptotic normality of  $f_n(x)$  will be proved together with that of the regression estimator  $m_n(x)$  of  $m(x)$  to be discussed later.

We first state a result from Real Analysis and a probability inequality to be used in the proof of [Theorem 9.2.1](#) before proving the theorem.

**Theorem 9.2.3.** *If  $K$  is of bounded variation on  $(-\infty, \infty)$  and  $\int |K(u)| du < \infty$ , then  $\lim_{u \rightarrow \pm\infty} K(u) = 0$ .*

**Theorem 9.2.4** (Dvoretzky-Kiefer-Wolfowitz). *If  $F_n$  is the empirical cdf of a random sample of size  $n$  from  $F$ , then there exists a constant  $C$  so that*

$$P\left[\sup_x |F_n(x) - F(x)| > a\right] \leq C \exp(-2na^2).$$

*Proof of Theorem 9.2.1.* Note that  $\sup_x |f_n(x) - f(x)| \leq A_n + B_n$ , where

$$\begin{aligned} B_n &= \sup_x |\mathbb{E}[f_n(x)] - f(x)| = \sup_x \left| \int_{-\infty}^{\infty} K(u) \{f(x - h_n u) - f(x)\} du \right| \\ &= \sup_x \left| \int_{-\infty}^{\infty} \left[ -h_n u f'(x) + (1/2) h_n^2 u^2 f''(x - \lambda h_n u) \right] K(u) du \right| \\ &\leq (1/2) h_n^2 \sigma_K^2 \sup_x |f''(x)| = O(h_n^2) = o(1), \end{aligned}$$

using Assumptions 1(a, b) and 2, and

$$\begin{aligned} A_n &= \sup_x |f_n(x) - \mathbb{E}[f_n(x)]| \\ &= \sup_x \left| (nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n) - \mathbb{E} \left[ (nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n) \right] \right| \\ &= h_n^{-1} \sup_x \left| \int_{-\infty}^{\infty} K((x - y)/h_n) d[F_n(y) - F(y)] \right| \\ &= h_n^{-1} \sup_x \left| \int_{-\infty}^{\infty} -[F_n(y) - F(y)] K((x - y)/h_n) dK((x - y)/h_n) \right| \\ &\leq (\mu/h_n) \sup_x |F_n(x) - F(x)|, \end{aligned}$$

where  $\mu = \int_{-\infty}^{\infty} K$  is the total variation of  $K$  on  $(-\infty, \infty)$ , having used integration by parts and observing that  $\lim_{y \rightarrow \pm\infty} h_n^{-1} K((x - y)/h_n) \{F_n(y) - F(y)\} = 0$ , all this by using Theorem 9.2.3. We shall now show that for arbitrary  $\varepsilon > 0$ ,  $\sum_{n=1}^{\infty} P[\sup_x |f_n(x) - \mathbb{E}[f_n(x)]| > \varepsilon] < \infty$  and then  $A_n \rightarrow 0$  a.s., will follow by the Borel-Cantelli Lemma. To this end, we use Theorem 9.2.4 to see that

$$\begin{aligned} P\left[\sup_x |f_n(x) - \mathbb{E}[f_n(x)]| > \varepsilon\right] &\leq P\left[(\mu/h_n) \sup_x |F_n(x) - F(x)| > \varepsilon\right] \\ &= P\left[\sup_x |F_n(x) - F(x)| > h_n \varepsilon / \mu\right] \\ &\leq C \exp[-2n(h_n \varepsilon / \mu)^2] = C \exp[-2(\varepsilon/\mu)^2 nh_n^2]. \end{aligned}$$

Since  $nh_n^2/\log n \rightarrow \infty$  as  $n \rightarrow \infty$ ,  $(\varepsilon/\mu)^2 nh_n^2 > \log n$  (ie,  $\exp[-2(\varepsilon/\mu)^2 nh_n^2] < 1/n^2$  for large  $n$ ) and  $\sum_{n=1}^{\infty} 1/n^2 < \infty$ .  $\square$

Functions of bounded variation, Stieltjes integral and integration by parts are discussed in Appendix A.3.

*Remark 9.2.1.* Strong uniform consistency of  $f_n$  holds under milder conditions than stated above. It is enough to assume:

- 1\*.  $K$  is a pdf of bounded variation.
- 2\*.  $f$  is uniformly continuous (and is therefore bounded).
- 3\*.  $h_n \downarrow 0$  and  $nh_n^2/\log n \rightarrow \infty$  as  $n \rightarrow \infty$ .

To prove [Theorem 9.2.1](#) under these conditions, note that only 1\* and 3\* are used in proving  $A_n \rightarrow 0$  a.s. in the above proof, so we only need to modify the proof of the fact that  $B_n = \sup_x |E[f_n(x)] - f(x)| \rightarrow 0$  as  $n \rightarrow \infty$  using  $\int_{-\infty}^{\infty} K(u) du < \infty$  and 2\*. The main points to note are that under these conditions, for every  $\varepsilon > 0$ ,

- (i) we can choose  $M$  such that  $\int_{|u|>M} K(u) < \varepsilon / \{4 \sup_x f(x)\}$ , so that

$$\int_{|u|>M} |f(x - h_n u) - f(x)| K(u) du < \varepsilon/2 \quad \text{for all } x, \text{ and}$$

- (ii) we can choose  $\delta_\varepsilon > 0$  such that

$$|h| < \delta_\varepsilon \implies \sup_x |f(x - h) - f(x)| < \varepsilon/2.$$

The details of the proof is left as an *exercise*.

*Remark 9.2.2.* For rates of convergence of kernel estimators of density functions and their derivatives, see Bhattacharya [48]. An estimator of Fisher-information of a location family of unknown form has also been constructed in this paper.

### 9.3 Regression Estimation

The regression function  $Y$  on  $X$  is  $m(\cdot)$ , where

$$m(x) = \varphi(x)/f_X(x), \quad \varphi(x) = \int_{-\infty}^{\infty} y f_{XY}(x, y) dy.$$

Since  $f_X(x)$  is estimated by  $f_n(x)$  given in Eq. (2), the main thing is to estimate  $\varphi(x)$ . Analogous to  $f(x) = F'(x)$ ,  $\varphi(x)$  can be expressed as

$$\begin{aligned} \varphi(x) &= \frac{d}{dx} \int_{t=-\infty}^x \left[ \int_{y=-\infty}^{\infty} y f_{XY}(t, y) dy \right] dt \\ &= \lim_{h \downarrow 0} h^{-1} \int_{x-h/2}^{x+h/2} \int_{y=-\infty}^{\infty} y f_{XY}(t, y) dy dt \\ &= \lim_{h \downarrow 0} h^{-1} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{[x-h/2, x+h/2]}(t) y f_{XY}(t, y) dy dt \\ &= \lim_{h \downarrow 0} E \left[ h^{-1} I_{[x-h/2, x+h/2]}(X) Y \right] = \lim_{h \downarrow 0} E \left[ h^{-1} I_{[-1/2, 1/2]}((x-X)/h) Y \right], \end{aligned}$$

so a natural estimator of  $\varphi(x)$  is

$$\varphi_n(x) = (nh_n)^{-1} \sum_{i=1}^n I_{[-1/2, 1/2]}((x - X_i)/h_n) Y_i$$

with small  $h_n > 0$ , or more generally,

$$\varphi_n(x) = (nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n) Y_i, \quad (7)$$

using a pdf  $K$  as the kernel.

This leads to the kernel regression estimator

$$m_n(x) = \frac{(nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n) Y_i}{(nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n)} = \frac{\varphi_n(x)}{f_n(x)}, \quad (8)$$

where  $h_n \downarrow 0$  as  $n \rightarrow \infty$ .

The estimator  $m_n(\cdot)$  also has the strong uniform convergence property, which we state below without proof.

**Theorem 9.3.1.** *Suppose that the following conditions hold:*

1. **(a)**  $P[a \leq X \leq b, c \leq Y \leq d] = 1$  for some  $a < b$  and  $c < d$ ,
- (b)**  $f_X(x)$  is bounded away from 0 on  $[a, b]$ ,
- (c)**  $|f_{XY}(x_1, y) - f_{XY}(x_2, y)| \leq M|x_1 - x_2|$  for some  $M$  and for all  $x_1, x_2 \in [a, b]$  and  $y \in [c, d]$ .
2.  $K$  is a bounded symmetric pdf on  $[-1, 1]$ .
3.  $h_n \rightarrow 0$  and  $nh_n/\log n \rightarrow \infty$  as  $n \rightarrow \infty$ .

Then  $\sup_{a \leq x \leq b} |m_n(x) - m(x)| \rightarrow 0$  a.s. at the rate of  $r_n = h_n + \sqrt{\log n/(nh_n)}$  (ie,  $r_n^{-1} \sup_{a \leq x \leq b} |m_n(x) - m(x)|$  is bounded with probability 1).

We now state a theorem on the asymptotic normality of the bivariate sequence  $(f_n(x), \varphi_n(x))$  given by Eqs. (2) and (7). In particular, this will establish the asymptotic normality of  $f_n(x)$ , proving [Theorem 9.2.2](#), and the asymptotic normality of  $m_n(x) = \varphi_n(x)/f_n(x)$  given in Eq. (8), will follow by the delta method.

We shall make the following assumptions:

1. The second derivative  $f''_X$ ,  $m''$ , and  $v''$ , where  $v(x) = \text{Var}[Y|X = x]$  exist, and are bounded and continuous.
2. The kernel  $K$  is a symmetric pdf with  $\sigma_K^2 = \int u^2 K(u) du < \infty$  and  $\|K\|^2 = \int K^2(u) du < \infty$ .

**Theorem 9.3.2.** *Under Assumptions 1 and 2, with  $h_n = n^{-1/5} t$ ,*

$$\begin{aligned} & n^{2/5} \begin{bmatrix} f_n(x) - f(x) \\ \varphi_n(x) - \varphi(x) \end{bmatrix} \\ & \xrightarrow{\mathcal{L}} N_2 \left( \frac{1}{2} t^2 \sigma_K^2 \begin{bmatrix} f''(x) \\ \varphi''(x) \end{bmatrix}, t^{-1} \|K\|^2 f(x) \begin{bmatrix} 1 & m(x) \\ m(x) & v(x) + m^2(x) \end{bmatrix} \right). \end{aligned}$$

**Corollary 9.3.1.** *Convergence of the first coordinate gives us Theorem 9.2.2, and using the delta method with  $g(u, w) = w/u$ , we have*

$$\begin{aligned} n^{2/5}[m_n(x) - m(x)] &\xrightarrow{\mathcal{L}} N(\beta(x), \Psi(x)), \text{ where} \\ \beta(x) &= (1/2)t^2\sigma_K^2 \frac{\varphi''(x) - f''(x)m(x)}{f(x)} \text{ and } \Psi(x) = t^{-1}\|K\|^2 \frac{v(x)}{f(x)}. \end{aligned}$$

Verification of the formulas for  $\beta(x)$  and  $\Psi(x)$  is left as an exercise.

*Proof of Theorem 9.3.2.* We shall use the Cramér-Wold device, namely,

$$a_1 U_n + a_2 W_n \xrightarrow{\mathcal{L}} a_1 U + a_2 W \text{ for all } (a_1, a_2) \Rightarrow (U_n, W_n) \xrightarrow{\mathcal{L}} (U, W).$$

Let us, therefore, consider

$$\begin{aligned} \xi_{ni} &= a_1 h_n^{-1} K((x - X_i)/h_n) + a_2 h_n^{-1} K((x - X_i)/h_n) Y_i \\ &= h_n^{-1} K((x - X_i)/h_n)(a_1 + a_2 Y_i). \end{aligned}$$

Then for each  $n$ ,  $\{\xi_{ni}, i = 1, \dots, n\}$  are iid. Let

$$\mu_{ni} = E[\xi_{ni}], \sigma_{ni}^2 = \text{Var}[\xi_{ni}], A_n = \sum_{i=1}^n \mu_{ni} = n\mu_{n1}, \text{ and } B_n^2 = \sum_{i=1}^n \sigma_{ni}^2 = n\sigma_{n1}^2.$$

Letting  $m(x) = E[Y|X = x]$  and  $v(x) = \text{Var}[Y|X = x]$ , we have

$$\begin{aligned} \mu_{ni} &= \int \int h_n^{-1} K((x - t)/h_n)(a_1 + a_2 y) f_{Y|X}(y|t) f_X(t) dy dt \\ &= \int \int K(u)(a_1 + a_2 y) f_{Y|X}(y|x - h_n u) f_X(x - h_n u) dy du \\ &= \int K(u)\{a_1 + a_2 m(x - h_n u)\} f_X(x - h_n u) du, \end{aligned}$$

and similarly,

$$\sigma_{ni}^2 = h_n^{-1} \int K^2(u) [\{a_1 + a_2 m(x - h_n u)\}^2 + a_2^2 v(x - h_n u)] f(x - h_n u) du - \mu_{ni}^2.$$

Assumptions 1 and 2 allow us to expand  $f(x - h_n u)$ ,  $m(x - h_n u)$ , and  $v(x - h_n u)$  to second-order terms about  $x$  in the above expressions. After algebraic manipulations, this leads to

$$\begin{aligned} \mu_{ni} &= \{a_1 + a_2 m(x)\} f(x) + (1/2)h_n^2 [\alpha(x) + o(1)] \text{ and} \\ \sigma_{ni}^2 &= h_n^{-1} [\gamma(x) + o(1)], B_n = n\sigma_{n1}^2 = nh_n^{-1} [\gamma(x) + o(1)] \end{aligned} \tag{9}$$

where

$$\begin{aligned} \alpha(x) &= \sigma_K^2 [a_1 f''(x) + a_2 \{f(x)m''(x) + f''(x)m(x) + 2f'(x)m'(x)\}], \\ \gamma(x) &= \|K\|^2 [\{a_1 + a_2 m(x)\}^2 + a_2^2 v(x)] f(x). \end{aligned} \tag{9a}$$

To establish asymptotic normality of  $\sum_{i=1}^n (\xi_{ni} - \mu_{ni})/B_n$ , we now check the Lindeberg Condition:

$$\lim_{n \rightarrow \infty} B_n^{-2} \sum_{i=1}^n E \left[ I_{(\varepsilon B_n, \infty)} (|\xi_{ni} - \mu_{ni}|) (\xi_{ni} - \mu_{ni})^2 \right] = 0,$$

the verification of which is left as an *exercise*.

We thus have  $\sum_{i=1}^n (\xi_{ni} - \mu_{ni})/B_n \xrightarrow{\mathcal{L}} N(0, 1)$ .

Using Eqs. (9) and (9a),  $\sum_{i=1}^n (\xi_{ni} - \mu_{ni})/B_n$  simplifies to

$$\begin{aligned} & \sqrt{\frac{nh_n}{\gamma(x) + o(1)}} \left[ a_1 \left\{ f_n(x) - f(x) - (1/2)h_n^2 \sigma_K^2 f''(x)(1 + o(1)) \right\} \right. \\ & \quad \left. + a_2 \left\{ \varphi_n(x) - \varphi(x) - (1/2)h_n^2 \sigma_K^2 \varphi''(x)(1 + o(1)) \right\} \right], \end{aligned}$$

which  $\xrightarrow{\mathcal{L}} N(0, 1)$  for all  $(a_1, a_2)$ .

Writing

$$\gamma(x) = (a_1 \ a_2) \begin{bmatrix} 1 & m(x) \\ m(x) & v(x) + m^2(x) \end{bmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \|K\|^2 f(x),$$

this implies

$$\begin{aligned} & \begin{bmatrix} \sqrt{nh_n} \{f_n(x) - f(x)\} - (1/2) \sqrt{nh_n^5} \sigma_K^2 f''(x)(1 + o(1)) \\ \sqrt{nh_n} \{\varphi_n(x) - \varphi(x)\} - (1/2) \sqrt{nh_n^5} \sigma_K^2 \varphi''(x)(1 + o(1)) \end{bmatrix} \\ & \xrightarrow{\mathcal{L}} N_2 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \|K\|^2 f(x) \begin{bmatrix} 1 & m(x) \\ m(x) & v(x) + m^2(x) \end{bmatrix} \right). \end{aligned}$$

Taking  $h_n = n^{-1/5}t$ , we arrive at the desired result. □

## 9.4 Nearest Neighbor Approach

Consider the uniform kernel with bandwidth  $h_n$ . In both density estimation and regression estimation, the methods discussed will suffer from the drawback that for those  $x$  where  $f(x)$  is small, very few datapoints will have their  $X$ -values falling in the window  $x \pm h_n/2$ . As a result, very few datapoints will contribute to the construction of  $f_n(x)$  or  $m_n(x)$  where  $f(x)$  is small.

To overcome this difficulty, one may decide to let the data determine the bandwidth  $h_n$ . The idea is to choose  $h_n$  large enough so that *exactly*  $k = k_n$  of the  $X_i$ 's are included in the window  $x \pm h_n/2$ . Clearly, this would put the burden of MSE more on the bias when  $f(x)$  is small (causing  $h_n$  to be large) and more on the variance when  $f(x)$  is large (causing  $h_n$  to be small).

Let  $d_n = \inf\{h: \sum_{i=1}^n I_{[x-h, x+h]}(X_i) \geq k_n\}$ . Then  $\sum_{i=1}^n I_{[x-d_n, x+d_n]}(X_i) = k_n$ . Now if we let  $d_n$  play the role of  $h_n/2$  in the kernel estimation procedure with uniform kernel, then

$$\begin{aligned} f_n(x) &= (nh_n)^{-1} \sum_{i=1}^n I_{[x-h_n/2, x+h_n/2]}(X_i) \\ &= \frac{1}{n(2d_n)} \sum_{i=1}^n I_{[x-d_n, x+d_n]}(X_i) = \frac{k_n}{2nd_n}. \end{aligned}$$

This  $d_n$  is called the  $k_n$ -nearest neighbor ( $k_n$ -NN) distance from  $x$  and  $f_n(x) = k_n/(2nd_n)$  is called the  $k_n$ -NN estimator of  $f(x)$ .

Here  $k_n$  is the smoothing parameter and the optimal rate at which  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$  is  $O(n^{4/5})$  (ie, we choose  $k_n = [n^{4/5}t]$ ), where  $[\cdot]$  is the greatest integer function. This is under the second-order smoothness condition assumed earlier.

Analogous argument in regression estimation leads to the  $k_n$ -NN regression estimator

$$\begin{aligned} m_n(x) &= \frac{(n2d_n)^{-1} \sum_{i=1}^n I_{[x-d_n, x+d_n]}(X_i) Y_i}{(n2d_n)^{-1} \sum_{i=1}^n I_{[x-d_n, x+d_n]}(X_i)} \\ &= k_n^{-1} \sum_{i=1}^n I_{[x-d_n, x+d_n]}(X_i) Y_i = k_n^{-1} \sum_{\{i: |X_i - x| \leq d_n\}} Y_i, \end{aligned}$$

that is,  $m_n(x)$  is the mean of those  $k_n$  values of  $Y_i$  corresponding to the  $X_i$ 's which are  $k_n$  nearest neighbors of  $x$ .

Another way to express this is to

- (i) replace  $(X_i, Y_i)$  by  $(|X_i - x|, Y_i) := (Z_i, Y_i)$ ,  $i = 1, \dots, n$ ,
- (ii) rank  $Z_i = |X_i - x|$ ,  $i = 1, \dots, n$  as  $0 < Z_{n:1} < \dots < Z_{n:n}$  (strict inequality w.p. 1),
- (iii) let  $Y_{n:i}$  be the  $Y$ -value associated with  $Z_{n:i}$ , ie,  $Y_{n:i} = Y_j \iff Z_{n:i} = Z_j$ .

Then the  $k_n$ -NN estimator of  $m(x)$  is  $m_n(x) = k_n^{-1} \sum_{i=1}^{k_n} Y_{n:i}$ .

## 9.5 Curve Estimation in Higher Dimension

Let  $\mathbf{X}$  be a  $d$ -dim rv with pdf  $f$  and  $Y$  be a real-valued rv whose regression on  $\mathbf{X}$  is  $m(\mathbf{x}) = E[Y|\mathbf{X} = \mathbf{x}]$ . So far, we have discussed the problems of estimating the pdf  $f$  from iid observations on  $X$  and of estimating the regression function from iid observation  $(X_1, Y_1), \dots, (X_n, Y_n)$  when  $d = 1$ . We now consider these problems for  $d \geq 2$ . For this, the kernel and NN methods and the theoretical results for their asymptotics described so far, extend in a straightforward manner to higher dimensions. However, the actual sample size  $n$  needed for these estimators to perform reasonably well, increases rapidly with  $d$ .

For  $d$ -dim  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , we work with  $d$ -dim kernel  $K_d$  which is symmetric in each coordinate and define

$$f_n(\mathbf{x}) = \left( nh_n^d \right)^{-1} \sum_{i=1}^n K_d((\mathbf{x} - \mathbf{X}_i)/h_n) \text{ and}$$

$$m_n(\mathbf{x}) = \left( nh_n^d \right)^{-1} \sum_{i=1}^n K_d((\mathbf{x} - \mathbf{X}_i)/h_n) Y_i / f_n(\mathbf{x}).$$

In particular, we may choose  $K_d(u_1, \dots, u_d) = \prod_{j=1}^d K(u_j)$  where  $K$  is a pdf on  $\mathbb{R}$  which is symmetric about 0.

To understand the difficulty of high dimensionality, let  $K_d$  be the uniform kernel on  $[-1/2, 1/2]^d$ . Then the only datapoints contributing to the construction of  $f_n(\mathbf{x})$  and  $m_n(\mathbf{x})$  are those with  $X_i - \mathbf{x} \in [-h_n/2, h_n/2]^d$ , the number of which is of the order of  $nh_n^d$ . Without going into detailed calculations, it is easy to see that the bias of  $f_n(\mathbf{x})$  or  $m_n(\mathbf{x})$  will still be of the order of  $h_n^2$  (as in the case of  $d = 1$ ), but their variances will be of the order of  $1/(nh_n^d)$ . This will result in mean square error of the form

$$A(\mathbf{x})h_n^4 + B(\mathbf{x})/(nh_n^d) + R_n(\mathbf{x}).$$

Neglecting the remainder term, we see that

$$n^{4/(d+4)}MSE \approx A(\mathbf{x})(n^{1/(d+4)}h_n)^4 + B(\mathbf{x})(n^{1/(d+4)}h_n)^{-d},$$

which blows up if  $n^{1/(d+4)}h_n$  either  $\rightarrow 0$  or  $\infty$ . Thus the optimal  $h_n = n^{-1/(d+4)}t$  with a suitable  $t$ , resulting in  $MSE = O(n^{-4/(d+4)})$ . This shows how the rate at which the MSE converges to 0 slows down as  $d$  increases. In the Curve Estimation literature, this phenomenon is called “Curse of Dimensionality.”

For example, if the  $X_i$ 's are sampled from  $N_d(\mathbf{0}, \mathbf{I})$ , then in order to estimate this density at  $\mathbf{x} = \mathbf{0}$  by the kernel method using normal kernel and optimal bandwidth, the sample size  $n$  needed for

$$\text{Relative } MSE[f_n(\mathbf{x})] = E[\{f_n(\mathbf{x}) - f(\mathbf{x})\}^2]/f^2(\mathbf{x})$$

to be  $< 0.1$  increases from  $n = 4$  for  $d = 1$  to  $n = 67$  for  $d = 3$  to  $n = 10,700$  for  $d = 7$  (see [49]).

## 9.6 Curve Estimation Using Local Polynomials

Let us first briefly review the result on the kernel estimator  $m_n(x)$  of the unknown regression function  $m(x)$  as given in Corollary 9.3.1, where asymptotic normality of the estimator  $m_n(x)$  is established after suitable renormalization. Asymptotic bias is

$$\text{asymptotic bias: } \left( h_n^2/2 \right) \sigma_K^2 \frac{\varphi''(x) - f''(x)m(x)}{f(x)},$$

where  $h_n = n^{-1/5}t$ ,  $t > 0$ . Since  $\varphi(x) = m(x)f(x)$ , we get  $\varphi''(x) - f''(x)m(x) = m''(x)f(x) + 2m'(x)f'(x)$ . From [Corollary 9.3.1](#) we get the asymptotic bias and asymptotic variance of  $m_n(x)$  as

$$\text{asymptotic bias: } \left(h_n^2/2\right)\sigma_K^2[m''(x) + 2m'(x)f'(x)/f(x)], \text{ and}$$

$$\text{asymptotic variance: } (nh_n)^{-1}\|K\|^2\nu(x)/f(x),$$

where  $\nu(x) = \text{Var}[Y|X = x]$ . Note that the bias term involves  $f'(x)$  (ie, it depends on the smoothness of the marginal density of  $X$ ). In other words, the design of the  $X_i$ 's enters the picture in regression estimation. Moreover, the estimator  $m_n(x)$  has another aspect which makes it difficult to use in practice. Suppose that  $f$  is supported on a compact interval which we take to be  $[0, 1]$  without loss of generality. The bias of the estimator  $m_n(x)$  at or near the boundary points 0 or 1 is of order  $h_n$  and not  $h_n^2$ . Thus the regression estimate is less reliable at or near the boundary points. The local polynomial method seeks to remove these negative aspects of the kernel estimator  $m_n$ .

In order to simplify notations let us denote

$$w_i(x) = \frac{K((x - X_i)/h_n)}{\sum_{j=1}^n K((x - X_j)/h_n)}.$$

Thus we may rewrite the kernel estimator as  $m_n(x) = \sum_{i=1}^n w_i(x)Y_i$ . Since  $w_i(x) \geq 0$  for all  $i$ , and  $\sum_{i=1}^n w_i(x) = 1$ ,  $m_n(x)$  is a weighted average of  $Y_i$ 's with weights  $w_i(x)$ . If the kernel  $K$  is compactly supported, say on  $[-1/2, 1/2]$ , then  $w_i(x) = 0$  whenever  $|X_i - x| > h_n/2$ . Note that a regression model for  $(X_i, Y_i)$  is of the form  $Y_i = m(X_i) + \varepsilon_i$ , where  $\{\varepsilon_i\}$  are independent with  $E[\varepsilon_i|X_i] = 0$  for all  $i$ . If  $X_i$  is in the neighborhood  $N_n(x) = \{u: |u - x| \leq h_n/2\}$  (ie,  $|X_i - x| \leq h_n/2$ ) and we approximate  $m(X_i)$  by a constant  $\beta_0 = m(x)$ , then the regression model takes the approximate form  $Y_i = \beta_0 + \varepsilon_i$  when  $X_i$  is in  $N_n(x)$ . If we obtain an estimate of  $\beta_0$  from a weighted least squares criterion of the form  $\sum_{i=1}^n (Y_i - \beta_0)^2 w_i(x)$ , which is minimized with respect to  $\beta_0$ , then  $m_n(x)$  is that value of  $\beta_0$  which minimizes this local weighted least squares.

Now if we approximate  $m(X_i)$  by a straight line instead of a constant when  $X_i$  is in  $N_n(x)$ , then a simple Taylor series expansion yields

$$m(X_i) = \beta_0 + \beta_1(X_i - x) + O((X_i - x)^2),$$

where  $\beta_0 = m(x)$  and  $\beta_1 = m'(x)$ . Thus the regression model described in the last paragraph can be approximately expressed as  $Y_i = \beta_0 + \beta_1(X_i - x) + \varepsilon_i$  when  $X_i$  is in  $N_n(x)$ . The parameters  $\beta_0$  and  $\beta_1$  can be estimated by using the method of weighted least squares. If the estimate of  $\beta_0$  is  $\hat{\beta}_0$ , then our estimate of  $m(x)$  is  $\hat{\beta}_0$ . More formally, we seek to minimize  $Q = \sum_{i=1}^n [Y_i - \beta_0 - \beta_1(X_i - x)]^2 w_i(x)$  with respect to  $\beta_0$  and  $\beta_1$ . If we differentiate  $Q$  with respect to  $\beta_0$  and  $\beta_1$  and equate the derivatives to 0, we are led to the equations

$$\beta_0 + \sum w_i(x)(X_i - x)\beta_1 = \sum w_i(x)Y_i,$$

$$\sum w_i(x)(X_i - x)\beta_0 + \sum w_i(x)(X_i - x)^2\beta_1 = \sum w_i(x)(X_i - x)Y_i.$$

We can get an explicit expression for the estimate of  $\beta_0$  from these equations after some tedious algebra. Let us simplify some notations by denoting

$$c_s(x) = \sum w_i(x)(X_i - x)^s, \quad s = 1, 2, 3.$$

Then the estimate of  $\beta_0$  is

$$\hat{\beta}_0 = \frac{\sum w_i(x)[c_2(x) - c_1(x)(X_i - x)]Y_i}{\sum w_i(x)[c_2(x) - c_1(x)(X_i - x)]}.$$

Thus the local linear estimate of  $m(x)$  is given by

$$m_n^{(LL)}(x) = \frac{\sum w_i(x)[c_2(x) - c_1(x)(X_i - x)]Y_i}{\sum w_i(x)[c_2(x) - c_1(x)(X_i - x)]} = \sum l_i(x)Y_i, \text{ where}$$

$$l_i(x) = \frac{w_i(x)[c_2(x) - c_1(x)(X_i - x)]}{\sum w_j(x)[c_2(x) - c_1(x)(X_j - x)]} = \frac{w_i(x)[c_2(x) - c_1(x)(X_i - x)]}{c_2(x) - c_1(x)^2}.$$

The local linear estimate of  $m(x)$  is also a linear combination of  $Y_i$ 's with weights  $l_i(x)$ , and the weights sum to 1 (ie,  $\sum l_i(x) = 1$ ). These weights also have the property  $\sum l_i(x)(X_i - x) = 0$ . This is unlike the regular kernel estimate  $m_n(x) = \sum w_i(x)Y_i$  where the weights  $w_i(x)$  do not necessarily satisfy the equation  $\sum w_i(x)(X_i - x) = 0$ . Why is this property important? It provides a correction term in the bias part. In order to see this, let us write

$$m_n^{(LL)}(x) - m(x) = \sum l_i(x)[Y_i - m(X_i)] + \sum l_i(x)[m(X_i) - m(x)].$$

The first sum in the last expression contributes toward the asymptotic variance of  $m_n^{(LL)}(x)$ , whereas the second sum contributes toward the asymptotic bias. We now examine the second sum in some detail. Let us assume that  $K$  is supported on a compact interval, say  $[-1/2, 1/2]$ , and it is symmetric about 0,  $m$  is twice differentiable,  $m''$  is continuous on the compact interval  $[a, b]$  on which  $m$  is being estimated and  $f$  is continuous on  $[a, b]$ . Under these conditions it can be shown that  $\inf_i l_i(x) \geq 0$  with probability converging to 1. Since  $m''$  is continuous on a compact interval, it is also uniformly continuous. Since  $w_i(x) = 0$  whenever  $|X_i - x| > h_n/2$ , using a two-term Taylor expansion of  $m(X_i)$  around  $m(x)$ , we have

$$\begin{aligned} \sum l_i(x)[m(X_i) - m(x)] &= \sum l_i(x)[(X_i - x)m'(x) + (1/2)(X_i - x)^2m''(x)] + o_P(h_n^2) \\ &= (1/2) \sum l_i(x)(X_i - x)^2m''(x) + o_P(h_n^2), \end{aligned}$$

because the sum involving the linear term with  $X_i - x$  is exactly equal to zero (ie,  $\sum l_i(x)(X_i - x) = 0$ ). A similar argument for the bias part for  $m_n(x)$  (ie,  $\sum w_i(x)[m(X_i) - m(x)]$ ) would involve a term of the form  $\sum w_i(x)(X_i - x)$ . In order for this term to be of order  $h_n^2$ , it is required that  $f$  is differentiable and  $f'$  is continuous. Such a condition is not needed for the local linear estimate.

Let us now examine the bias term for the local linear estimator  $m_n^{(LL)}(x)$  a bit more closely and establish that  $\sum l_i(x)(X_i - x)^2 = h_n^2[\sigma_K^2 + o_P(1)]$ . Then the bias term would be  $(1/2)h_n^2\sigma_K^2m''(x) + o_P(h_n^2)$ .

Using the expression for  $l_i(x)$  we get

$$\begin{aligned}\sum l_i(x)(X_i - x)^2 &= \sum \frac{w_i(x)[c_2(x) - c_1(x)(X_i - x)](X_i - x)^2}{c_2(x) - c_1(x)^2} \\ &= \frac{c_2(x)^2 - c_1(x)c_3(x)}{c_2(x) - c_1(x)^2}.\end{aligned}$$

For any nonnegative integer  $r$ , using Chebychev's inequality, it can be established that

$$\begin{aligned}\left(nh_n^{s+1}\right)^{-1} \sum K((x - X_i)/h_n)(X_i - x)^s \\ = (nh_n)^{-1} \sum K((x - X_i)/h_n)\{(X_i - x)/h_n\}^s \xrightarrow{P} \int u^s K(u) du f(x).\end{aligned}$$

Noting that  $\int u^s K(u) du = 0$  for  $s = 1$  and  $s = 3$  (since  $K$  is symmetric about 0), we have

$$\begin{aligned}c_1(x) &= h_n \left[ \int u K(u) du + o_P(1) \right] = o_P(h_n), \\ c_2(x) &= h_n^2 \left[ \int u^2 K(u) du + o_P(1) \right] = h_n^2 [\sigma_K^2 + o_P(1)], \text{ and} \\ c_3(x) &= h_n^3 \left[ \int u^3 K(u) du + o_P(1) \right] = o_P(h_n^3).\end{aligned}$$

Plugging in these approximations for  $c_1(x)$ ,  $c_2(x)$ , and  $c_3(x)$  in the expression for  $\sum l_i(x)(X_i - x)^2$ , we have

$$\begin{aligned}\sum l_i(x)(X_i - x)^2 &= \frac{c_2(x)^2 - c_1(x)c_3(x)}{c_2(x) - c_1(x)^2} \\ &= \frac{h_n^4 [\sigma_K^2 + o_P(1)]^2 - o_P(h_n^4)}{h_n^2 [\sigma_K^2 + o_P(1)] - o_P(h_n^2)} = h_n^2 [\sigma_K^2 + o_P(1)].\end{aligned}$$

Using arguments similar to the ones in proving Corollary 9.3.1, we can also get a similar result for the local linear estimate  $m_n^{(LL)}(x)$ .

**Theorem 9.6.1.** *Assume that*

- (a) *on a compact interval  $[a, b]$ ,  $m$  is twice differentiable, and  $m''$ ,  $v$ , and  $f$  are continuous,*
- (b) *the kernel  $K$  is a bounded pdf on a compact interval  $[-c, c]$  and is symmetric about zero, and*
- (c)  $\sup_{x \in [a, b]} E[Y^4 | X = x] < \infty$ .

Taking  $h_n = n^{-1/5}t$ , it can be shown that for  $x \in [a, b]$ ,

$$\begin{aligned} n^{2/5} [m_n^{(LL)}(x) - m(x)] &\xrightarrow{\mathcal{L}} N(\beta^{(LL)}(x), \Psi^{(LL)}(x)), \text{ where} \\ \beta^{(LL)}(x) &= (1/2)t^2\sigma_K^2 m''(x) \text{ and } \Psi(x) = t^{-1}\|K\|^2 \frac{\nu(x)}{f(x)}. \end{aligned}$$

Note that the asymptotic variances of  $m_n(x)$  and  $m_n^{(LL)}(x)$  are the same, but their asymptotic biases are different.

*Remark 9.6.1.* In this section, a few aspects of local polynomial method have been highlighted. More applications and details can be found in the book by Fan and Gijbels [50].

1. Instead of the local linear method for estimating  $m(x)$ , one may consider local polynomial estimation by minimizing  $Q = \sum [Y_i - \beta_0 - \beta_1(X_i - x) - \dots - \beta_p(X_i - x)^p]^2 w_i(x)$  with respect to  $\beta_0, \beta_1, \dots, \beta_p$ . Then the local polynomial estimate of  $m(x)$  is  $m_n^{(LP)}(x) = \hat{\beta}_0$ .
2. If the independent variable  $X$  is vector valued, one can obtain a local linear (or more generally a local polynomial) estimate of  $m(x)$ . In this case, one minimizes  $Q = \sum [Y_i - \beta_0 - \beta_1^T(X_i - x)]^2 w_i(x)$  with respect to  $\beta_0$  and  $\beta_1$ , and as before the estimate of  $m(x)$  is given by  $\hat{\beta}_0$ .
3. It is also possible to carry out density estimation using the local linear or local polynomial method.
4. The choice of bandwidth  $h_n$  is crucial as in any other curve estimation problem. A cross-validation method can also be employed to obtain an appropriate value of  $h_n$ . The method of cross-validation is described in the next subsection in the context of kernel density and regression estimation problems.

### 9.6.1 Choice of $h_n$ in Density Estimation

Ideally, we would choose  $h_n$  which minimizes the Integrated Square-Error

$$\int (f_n - f)^2 = \int f_n^2 - 2 \int f_n f + \int f^2,$$

in which the last term does not depend on  $f_n$ . So the aim is to minimize  $R(f_n) = \int f_n^2 - 2 \int f_n f$ . The first term  $\int f_n^2$  is calculated directly from  $f_n$ . So the main thing is to estimate the second term from the data. The idea of “leave-one-out” is used for this purpose. Let

$$f_{n,-i}(x) = \left[ \{(n-1)h_n\}^{-1} \sum_{j \neq i=1}^n K((x - X_j)/h_n) \right] \quad \text{for each } i = 1, \dots, n. \quad (10)$$

Then

$$\begin{aligned} \mathbb{E}\left[n^{-1} \sum_{i=1}^n f_{n,-i}(X_i)\right] &= \mathbb{E}[f_{n,-n}(X_n)] = \mathbb{E}\mathbb{E}[f_{n,-n}(X_n)|X_1, \dots, X_{n-1}] \\ &= \mathbb{E} \int f_{n,-n}(x) f(x) dx = \mathbb{E} \int f_n(x) f(x) dx, \end{aligned}$$

because

$$\begin{aligned} \mathbb{E}[f_{n,-n}(x)] &= \mathbb{E}\left[\{(n-1)h_n\}^{-1} \sum_{j=1}^{n-1} K((x-X_j)/h_n)\right] \\ &= \mathbb{E}\left[(nh_n)^{-1} \sum_{j=1}^n K((x-X_j)/h_n)\right] = \mathbb{E}[f_n(x)]. \end{aligned}$$

Thus the estimator

$$\hat{R}(f_n) = \int f_n^2 - (2/n) \sum_{i=1}^n f_{n,-i}(X_i) \quad (11)$$

has the property:  $\mathbb{E}[\hat{R}(f_n)] = \mathbb{E}[R(f_n)]$ .

Hence minimizing  $\mathbb{E}[\hat{R}(f_n)]$  is equivalent to minimizing  $\mathbb{E}[R(f_n)]$ , so a choice of  $h_n$  which minimizes  $\hat{R}(f_n)$  will, hopefully, approximate the optimal  $h_n$  which minimizes  $R(f_n)$  itself.

We now describe the actual cross-validation procedure. First, for computational facility, we express

$$\begin{aligned} \int f_n^2 &= (nh_n)^{-2} \sum_{i=1}^n \sum_{j=1}^n \int K((X_i-x)/h_n) K((x-X_j)/h_n) dx \\ &= (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n \int K(X_i/h_n - u) K(u - X_j/h_n) du \\ &= (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n \int K((X_i - X_j)/h_n - (u - X_j/h_n)) K(u - X_j/h_n) du \\ &= (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n K^{(2)}((X_i - X_j)/h_n), \end{aligned}$$

where  $K^{(2)}$  is the convolution of  $K$  with itself, that is

$$K^{(2)}(z) = \int K(z-u) K(u) du = \text{pdf of } Z_1 + Z_2,$$

where  $Z_1$  and  $Z_2$  are iid with pdf  $K$ . Next rewrite the term  $n^{-1} \sum_{i=1}^n f_{n,-i}(X_i)$  with a minor modification to obtain

$$\begin{aligned} n^{-1} \sum_{i=1}^n f_{n,-i}(X_i) &= n^{-1} \sum_{i=1}^n [(n-1)h_n]^{-1} \sum_{j \neq i=1}^n K((X_i - X_j)/h_n) \\ &\approx (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n K((X_i - X_j)/h_n) - (nh_n)^{-1} K(0). \end{aligned}$$

Putting the two terms together in Eq. (11), we have

$$\begin{aligned} \hat{R}(f_n) &\approx \hat{R}_1(f_n) \\ &= (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n K^{(2)}((X_i - X_j)/h_n) \\ &\quad - 2 \left\{ (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n K((X_i - X_j)/h_n) - (nh_n)^{-1} K(0) \right\} \\ &= (n^2 h_n)^{-1} \sum_{i=1}^n \sum_{j=1}^n K^*((X_i - X_j)/h_n) + 2(nh_n)^{-1} K(0), \end{aligned}$$

where  $K^*(\cdot) = K^{(2)}(\cdot) - 2K(\cdot)$ .

The least squares cross-validation procedure is to choose the bandwidth as the value of  $h_n$  which minimizes  $\hat{R}_1(f_n)$  (see [49, 51]).

Let  $I(f_n; h_n) = \int (f_n - f)^2$ , where  $f_n$  is given by Eq. (2) and let  $h_n(CV)$ ,  $h_n(opt)$  denote, respectively, the  $h_n$  obtained by minimizing  $\hat{R}_1(f_n)$  and the unknown optimal  $h_n$  which minimizes  $I(f_n; h_n)$  for the given data. Then the following optimality property holds for  $h_n(CV)$  whenever  $f$  is bounded and  $K$  satisfies some mild conditions:

$$\lim_{n \rightarrow \infty} \frac{I(f_n; h_n(CV))}{I(f_n; h_n(opt))} = 1 \quad \text{with probability 1.}$$

The above discussion was for  $d = 1$ . All of this goes through for  $d \geq 2$ , by replacing  $h_n$  by  $h_n^d$  in the formula for  $\hat{R}_1(f_n)$ .

### 9.6.2 Regression Estimation

In regression estimation, analogous to density estimation, our goal may be to minimize

$$\begin{aligned} d(m_n; h_n) &= \int (m_n - m)^2 wf \\ &= \int m_n^2 wf - 2 \int m_n m wf + \int m^2 wf, \end{aligned}$$

where  $f$  is the pdf of  $X$  and  $w(x) > 0$  is a weight function. Again, the last terms does not depend on  $m_n$ , while the first two terms involve the unknown  $m$  and  $f$  and therefore, need to be estimated.

To motivate the proposed procedure, let  $(X_0, Y_0)$  be an observation on  $(X, Y)$  which is independent of the data  $(X_1, Y_1), \dots, (X_n, Y_n)$  and then rewrite the two terms as

$$\begin{aligned} \int m_n^2 wf - 2 \int m_n m wf &= E\left[\left\{m_n^2(X_0) - 2m_n(X_0)E(Y_0|X_0)\right\}w(X_0)|Data\right] \\ &= E\left[E\left(\left\{m_n^2(X_0) - 2m_n(X_0)Y_0\right\}w(X_0)|X_0, Data\right)|Data\right] \\ &= E\left[\left\{m_n^2(X_0) - 2m_n(X_0)Y_0\right\}w(X_0)|Data\right] \\ &= E\left[\{Y_0 - m_n(X_0)\}^2 w(X_0)|Data\right] - E\left[Y_0^2 w(X_0)\right]. \end{aligned}$$

Since  $E[Y_0^2 w(X_0)]$  does not depend on the choice of  $h_n$ , we only need to estimate  $E[\{Y_0 - m_n(X_0)\}^2 w(X_0)|Data]$  and then choose  $h_n$  to minimize this estimate.

Again using the “leave-one-out” method, we construct the estimate

$$\hat{d}_1(m_n; h_n) = n^{-1} \sum_{i=1}^n \{Y_i - m_{n,-i}(X_i)\}^2 w(X_i)$$

for  $E[\{Y_0 - m_n(X_0)\}^2 w(X_0)|Data]$ , where

$$m_{n,-i}(x) = [(n-1)h_n]^{-1} \sum_{j \neq i=1}^n K\left(\frac{(x-X_j)/h_n}{h_n}\right) Y_j / f_{n,-i}(x)$$

with  $f_{n,-i}(x)$  as in Eq. (10).

Then the cross-validated choice  $h_n(CV)$  and the optimal choice  $h_n(opt)$  of  $h_n$  are, respectively, the minimizer of  $\hat{d}_1(m_n; h_n)$  and  $d(m_n; h_n)$ . However, for technical reason, the choice of  $h_n$  is restricted to  $[C_1 n^{-\delta}, C_2]$  for some constants  $C_1, C_2$  and  $\delta > 0$ . This technical condition is also in conformity with the practice. To see this, assume that the kernel  $K$  is supported on  $[-1/2, 1/2]$ . Then the nonparametric regression estimate at  $x$  is simply a weighted average of  $Y_i$ 's for which  $X_i$  is in the interval  $[x - h_n/2, x + h_n/2]$ . Now, if  $h_n$  is too small and  $x$  is not one of the  $X_i$ 's, then this interval may not have  $X$ -observations, thus there are no observations to average on and one cannot get a nonparametric estimate of  $m$  at  $x$ .

Suppose that  $f$  is supported on a compact set in  $\mathbb{R}^d$ , on which it is bounded away from 0, and that  $f$  and  $m$  are continuous on this set. Then under assumption of boundedness of conditional moments of all orders of  $Y$  given  $X$ , the following property holds for  $h_n(CV)$ :

$$\frac{d(m_n; h_n(CV))}{d(m_n; h_n(opt))} \xrightarrow{P} 1.$$

## 9.7 Estimation of Survival Function and Hazard Rates Under Random Right-Censoring

Let  $(T_1, C_1), \dots, (T_n, C_n)$  be iid pairs of positive-valued rv's where for each  $i$ ,  $T_i$  and  $C_i$  are independent. For each  $i$ ,  $T_i$  is the *survival time* (ie, the time until death of a sample subject or the failure time of a sample equipment) and  $C_i$  is the *censoring time* (ie, time at which observation is stopped for this sample unit). Thus the observed data consist of  $(Y_1, \delta_1), \dots, (Y_n, \delta_n)$  where  $Y_i = T_i \wedge C_i$  and  $\delta_i = I(T_i \leq C_i)$ , so  $(Y_i, \delta_i) = (y, 0)$  means that the observation on the  $i$ th unit was censored at time  $y$  and  $(Y_i, \delta_i) = (y, 1)$  means that the observation continued until death or failure which occurred at time  $y$ . This is *random right censoring* of survival time. For simplicity, we shall assume that there are no ties.

Let  $F$  denote the common cdf of the  $T_i$ 's and let  $f$  denote the corresponding pdf.

**Definition 9.7.1.** The function  $S(t) = 1 - F(t) = P[T > t]$  is called the survival function and  $\lambda(t) = f(t)/[1 - F(t)]$  is called the hazard function or the hazard rate of the rv  $T$ .

Since  $\lambda(t) dt \approx P[t < T < t + dt | T > t]$ , it is also called the instantaneous failure rate.

The survival function  $S(t)$  and the hazard function  $\lambda(t)$  are related by the formula

$$S(t) = \exp\left[-\int_0^t \lambda(u) du\right], \quad (12)$$

the proof of which is left as an *exercise*.

### Estimation of the Survival Function

We now consider the problem of estimating the survival function  $S(t)$  from randomly right-censored data  $\{(Y_i, \delta_i), i = 1, \dots, n\}$ . Let  $Y_{n:1} < \dots < Y_{n:n}$  denote the order statistics of  $Y_1, \dots, Y_n$  and let  $\delta_{n:1}, \dots, \delta_{n:n}$  be defined by  $\delta_{n:i} = \delta_j \iff Y_{n:i} = Y_j$ . At each  $Y_{n:i}$ , either a death or a censoring occurs, depending on whether  $\delta_{n:i} = 1$  (death) or  $\delta_{n:i} = 0$  (censoring).

Consider the intervals  $I_i = (Y_{n:i-1}, Y_{n:i}]$ , taking  $Y_{n:0} = 0$ , and let

$\mathcal{R}(t) =$  Risk set at time  $t$  consisting of those who are still alive at time  $t-$ ,

$n_i = \#\mathcal{R}(Y_{n:i}) =$  number alive at time  $Y_{n:i}-$ ,

$d_i =$  number dying at time  $Y_{n:i} = \delta_{n:i}$ ,

$p_i = P[T > Y_{n:i} | T \geq Y_{n:i-1}]$  and  $q_i = 1 - p_i$ .

The natural estimates of  $q_i$  and  $p_i$  are

$$\hat{q}_i = d_i/n_i \text{ and } \hat{p}_i = 1 - \hat{q}_i = \begin{cases} 1 - 1/n_i & \text{if } \delta_{n:i} = 1 \\ 1 & \text{if } \delta_{n:i} = 0 \end{cases} = (1 - 1/n_i)^{\delta_{n:i}}.$$

The Product-Limit (PL) estimator of  $S(t)$  due to Kaplan and Meier [52] is

$$\hat{S}(t) = \prod_{i: Y_{n:i} \leq t} \hat{p}_i = \prod_{i: Y_{n:i} \leq t} (1 - 1/n_i)^{\delta_{n:i}}, \quad (13)$$

because  $n_i = \#\mathcal{R}(Y_{n:i-1}) = n - \#\{\text{dead or censored before } Y_{n:i-1}\} = n - i + 1$ .

*Remark 9.7.1.* On the interval  $(Y_{n:i-1}, Y_{n:i}]$ ,  $\hat{S}(t)$  remains unchanged if the observation at  $Y_{n:i}$  is censored and it is reduced by a factor of  $(n - i)/(n - i + 1)$  if  $Y_{n:i}$  is uncensored (death).

### Variance of $\hat{S}(t)$

Since  $n_i \hat{p}_i \sim \text{Bin}(n_i, p_i)$ ,

$$\text{Var}[\log \hat{p}_i] \approx \text{Var}[\hat{p}_i](d \log p_i / dp_i)^2 = (p_i q_i / n_i) \left(1/p_i^2\right) = q_i / (n_i p_i)$$

by the delta method. Hence

$$\text{Var}[\log \hat{S}(Y_{n:i})] = \text{Var}\left[\sum_{j=1}^i \log \hat{p}_j\right] = \sum_{j=1}^i \text{Var}[\log \hat{p}_j] \approx \sum_{j=1}^i \hat{q}_j / (n_j \hat{p}_j),$$

assuming  $\log \hat{p}_1, \log \hat{p}_2, \dots$  are independent. Using the delta method again, we have

$$\begin{aligned} \text{Var}[\hat{S}(Y_{n:i})] &= \text{Var}\left[\exp(\log \hat{S}_{n:i})\right] \approx \exp(2 \log \hat{S}_{n:i}) \text{Var}[\log \hat{S}_{n:i}] \\ &\approx \hat{S}^2(Y_{n:i}) \sum_{j=1}^i \hat{q}_j / (n_j \hat{p}_j). \end{aligned}$$

Thus

$$\text{Var}[\hat{S}(t)] \approx \hat{S}^2(t) \sum_{i: Y_{n:i} \leq t} \hat{q}_i / (n_i \hat{p}_i) = \hat{S}^2(t) \sum_{i: Y_{n:i} \leq t} \delta_{n:i} / \{n_i(n_i - 1)\}.$$

### Redistribute-to-the-Right Algorithm

This is another method, due to Efron [53], of calculating the PL estimator  $\hat{S}(t)$  given by Eq. (13). If we start with a sample of size  $n$ ,

- (i) first put a probability of  $1/n$  at each  $Y_{n:i}$ ;
- (ii) if  $Y_{n:i_1}$  is the first censored observation, redistribute the probability  $1/n$  assigned to  $Y_{n:i_1}$  equally to  $Y_{n:i_1+1}, \dots, Y_{n:n}$ , so that each of these observations now carries a probability of  $(1/n)(1 + 1/(n - i_1))$ ;
- (iii) if  $Y_{n:i_2}$  is the second censored observation, distribute the probability on  $Y_{n:i_2}$  equally to  $Y_{n:i_2+1}, \dots, Y_{n:n}$ , so that each of these observations now carries a probability of  $(1/n)(1 + 1/(n - i_1))(1 + 1/(n - i_2))$ , and so on.

A proof that this method leads to the same estimator as the one given by Eq. (13) is left as an *exercise*.

## Estimation of the Integrated Hazard Function

Estimation of the hazard function from censored data would involve estimation of the pdf of  $T$  from censored data, which is difficult. On the other hand, estimating the integrated hazard function  $\Lambda(t) = \int_0^t \lambda(u) du$  is straightforward, since  $S(t) = \exp[-\Lambda(t)]$  by Eq. (12). We therefore estimate  $\Lambda(t)$  by the estimator of  $-\log S(t)$ , that is,

$$\hat{\Lambda}_1(t) = -\log \hat{S}(t) = -\sum_{Y_{n:i} \leq y} \log(1 - \delta_{n:i}/(n - i + 1))$$

because  $(1 - 1/(n - i + 1))^{\delta_{n:i}} = 1 - \delta_{n:i}/(n - i + 1)$ . Also since  $-\log(1 - x) \approx x$  for small  $x$ , we have another estimator

$$\hat{\Lambda}_2(t) = \sum_{Y_{n:i} \leq t} \delta_{n:i}/(n - i + 1)$$

which is approximately the same as  $\hat{\Lambda}_1$ .

## Exercises

- 9.1. Show that a sufficient condition for a kernel estimator  $f_n(x)$  of  $f(x)$  with bandwidth  $h_n$  to be a consistent estimator is that  $h_n \downarrow 0$  and  $nh_n \rightarrow \infty$ .
- 9.2. Suppose that  $f$  is  $m$  times differentiable and  $f^{(m)}$  is bounded. Drop the condition that the kernel  $K$  is a pdf, but satisfies the conditions:  
 $\int K(u) du = 1$ ,  $\int u^r K(u) du = 0$ ,  $r = 1, \dots, m-1$ ,  $\int |u|^m K(u) du < \infty$ , and  
 $\int K^2(u) du < \infty$ .
  - (a) Find the asymptotic bias and variance of the estimator  
 $f_n(x) = (nh_n)^{-1} \sum_{i=1}^n K((x - X_i)/h_n)$  using such a kernel  $K$ .
  - (b) Determine the optimal rate at which  $h_n$  should tend to 0 and the corresponding rate of convergence of the MSE of  $f_n(x)$ .
- 9.3. Let  $m_n(x)$  be a kernel estimator of the regression function  $m(x)$  of  $Y$  on  $X$  at  $X = x$  based on a random sample of size  $n$ . Verify the formula for the mean  $\beta(x)$  and the variance  $\Psi(x)$  of the asymptotic distribution of  $n^{2/5}[m_n(x) - m(x)]$  given in the text.
- 9.4. Give a detailed proof of Theorem 9.2.1 under the milder conditions 1\*, 2\*, and 3\*.
- 9.5. Find the formulas of bias and variance of the  $k_n$ -NN estimators of a pdf  $f(x)$  and a regression function  $m(x)$ , and verify that the optimal rate at which  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$  is  $O(n^{4/5})$ .
- 9.6. Prove formula (12) on the relation between the survival function  $S(t)$  and the hazard function  $\lambda(t)$ .
- 9.7. Prove that the redistribute-to-the-right algorithm leads to the same estimator as the one given by Eq. (13).

# Statistical Functionals and Their Use in Robust Estimation

## 10.1 Introduction

Let  $\mathcal{F}_0$  be a family of cdf's in  $\mathbb{R}^d$ . Then  $T: \mathcal{F}_0 \rightarrow \mathbb{R}$  is called a statistical functional. Most statistical problems involve inference about such a  $T(F)$  for an unknown  $F \in \mathcal{F}_0$  based on  $T(F_n)$ , where  $F_n$  is the empirical distribution function of a random sample  $X_1, \dots, X_n$  from  $F$ . The behavior of  $T(F_n) - T(F)$  is, therefore, of interest. Study of statistical functionals was introduced by von Mises [54].

### Examples

In these examples, for simplicity,  $d = 1$ .

1.  $T(F) = E_F[g(X)] = \int g(x) dF(x)$  on  $\mathcal{F}_0 = \{F: E_F[|g(X)|] < \infty\}$ ,  $T(F_n) = n^{-1} \sum_{i=1}^n g(X_i)$ .
2.  $T(F) = E_F[(X - \xi(F))^k]$  on  $\mathcal{F}_0 = \{F: E_F|X|^k < \infty\}$ , where  $\xi(F) = E_F[X]$ ,  
 $T(F_n) = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^k$  where  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ .
3.  $T(F) = F^{-1}(p) = \inf\{x: F(x) \geq p\}$ ,  $0 < p < 1$ , on  $\mathcal{F}_0 = \{\text{all cdf's on } \mathbb{R}\}$ ,  $T(F_n) = X_{n: [np]}$ .
4.  $T(F) = E_F[g(X_1, \dots, X_r)]$  where  $g$  is symmetric in its coordinates,  
 $\mathcal{F}_0 = \{F: E_F[|g(X_1, \dots, X_r)|] < \infty\}$ ,  $T(F_n) = n^{-r} \sum_{i_1=1}^n \cdots \sum_{i_r=1}^n g(X_{i_1}, \dots, X_{i_r})$ .

In the last example,  $T(F_n)$  is called a  $V$ -statistic which differs from a  $U$ -statistic by its inclusion of all  $(i_1, \dots, i_r)$  rather than only those for which  $i_1 \neq \dots \neq i_r$  and then dividing the sum by  $n^r$  instead of  $n^{(r)} = n(n-1)\cdots(n-r+1)$ .

In Section 10.2, we shall introduce an expansion of  $\sqrt{n}[T(F_n) - T(F)]$  by means of “differentials” of  $T(F)$ , analogous to the expansion of  $\sqrt{n}[T(\hat{\theta}_n) - T(\theta)]$ . The leading term of this expansion will provide a functional delta method subject to the remainder term being  $o_P(1)$ . Postponing the issue of remainder terms, Sections 10.3–10.5 will be devoted to the  $L$ - and  $M$ -estimators. These estimators are called “robust” because their properties hold for a wide class of distributions, unlike estimators focused on squared-error loss or the maximum likelihood estimators which are susceptible to distributions with heavy tails. Finally, the issue of remainder terms is taken up in Section 10.6.

## 10.2 Functional Delta Method

We shall assume that  $\mathcal{F}_0$  is a collection of cdf's such that

- (i)  $F \in \mathcal{F}_0 \Rightarrow F_n \in \mathcal{F}_0$  and
- (ii)  $\mathcal{F}_0$  is convex (ie,  $F, G \in \mathcal{F}_0$  implies  $F + \lambda(G - F) \in \mathcal{F}_0$  for any  $\lambda \in [0, 1]$ ).

First consider a  $k$ -dim parameter  $\theta = \theta(F)$  being estimated by  $\hat{\theta}_n$  and a real-valued function  $T(\theta)$  of  $\theta$  being estimated by  $T(\hat{\theta}_n)$ . For large  $n$ , taking the leading term in the expansion of  $\sqrt{n}[T(\hat{\theta}_n) - T(\theta)]$ , we have

$$\sqrt{n}[T(\hat{\theta}_n) - T(\theta)] = \langle \nabla T(\theta), \sqrt{n}(\hat{\theta}_n - \theta) \rangle + R_n,$$

where  $R_n = o_P(1)$  if  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} W$ , a random vector, and  $T$  is continuously differentiable. In that case, the asymptotic distribution of  $\sqrt{n}[T(\hat{\theta}_n) - T(\theta)]$  is the same as that of the leading term in the expansion (ie,  $\sqrt{n}[T(\hat{\theta}_n) - T(\theta)] \xrightarrow{\mathcal{L}} \langle \nabla T(\theta), W \rangle$ ). This is commonly called the delta-method (see [Theorem 3.2.6, Chapter 3](#)). For a similar analysis of  $\sqrt{n}[T(F_n) - T(F)]$ , we need an analogous expansion which involves *differentiation of a functional  $T$  at  $F$* .

To understand the meaning of differentiation of  $T$  at  $F$ , let us again look at  $T: \mathbb{R}^k \rightarrow \mathbb{R}$  and examine the one-term Taylor expansion of  $T(\theta + \Delta) - T(\theta)$  for small  $\Delta$ :

$$T(\theta + \Delta) - T(\theta) = \sum_{r=1}^k T'(r; \theta) \Delta_r + o(\|\Delta\|), \quad (1)$$

where  $\|\Delta\| = (\sum_{r=1}^k \Delta_r^2)^{1/2}$  or  $\max_{1 \leq r \leq k} |\Delta_r|$ , and  $T'(r; \theta) = \partial T(u)/\partial u_r|_{u=\theta}$ .

The leading term is the differential  $L_\theta(\Delta) := \sum_{r=1}^k T'(r; \theta) \Delta_r$ , which is linear in  $\Delta$ .

Replacing  $\theta$  by  $F$ ,  $\theta + \Delta$  by  $G$ ,  $r$  by  $x$ ,  $\sum_{r=1}^k$  by  $\int_{x \in \mathbb{R}^d}$ , and  $\Delta_r$  by  $d\Delta(x) = d[G(x) - F(x)]$ , the expansion (1) would take the form

$$\begin{aligned} T(G) - T(F) &= L_F(G - F) + o(\|\Delta\|_\rho) \\ &= \int T'(x; F) d[G(x) - F(x)] + o(\|\Delta\|_\rho). \end{aligned}$$

Such an expansion would be valid if there exists a linear functional  $L_F$  on  $\mathcal{D} = \{c(G - F): c \in \mathbb{R} \text{ and } F, G \in \mathcal{F}_0\}$  which is not identically 0 and satisfies

$$\lim_{j \rightarrow \infty} \frac{T(G_j) - T(F) - L_F(G_j - F)}{\|G_j - F\|_\rho} = 0,$$

whenever  $\{G_j\}$  is a sequence in  $\mathcal{F}_0$  with  $\|G_j - F\|_\rho = \rho(G_j, F) \rightarrow 0$  as  $j \rightarrow \infty$ . Here  $\rho$  is a distance such that

$$\rho(F, F + t(G - F)) = |t|\rho(F, G) \iff \|t(G - F)\| = |t| \|G - F\|_\rho$$

holds for all  $t$  in  $\mathbb{R}$  and  $F, G \in \mathcal{F}_0$ .

In general, existence of  $L_F$  with this property depends on how the sequence  $\{G_j\}$  is allowed to approach  $F$  in the metric  $\rho$ . Three such schemes are described below. In each scheme, the linear functional  $L_F$ , if it exists, is called the *differential*.

- A.** *Gâteaux differentiability.*  $T$  is Gâteaux differentiable if the differential exists for every sequence  $G_j = F + t_j \Delta \in \mathcal{F}_0$ , where  $\Delta = G - F \in \mathcal{D}$  is fixed and  $t_j \rightarrow 0$ . Here  $\rho(G_j, F) = \|G_j - F\|_\rho = \|t_j \Delta\|_\rho = |t_j| \|\Delta\|_\rho \rightarrow 0$  irrespective of the metric  $\rho$ .
- B.**  *$\rho$ -Hadamard differentiability.*  $T$  is  $\rho$ -Hadamard differentiable if the differential exists for every sequence  $G_j = F + t_j \Delta_j \in \mathcal{F}_0$ , where  $\|\Delta_j - \Delta\|_\rho \rightarrow 0$  for a fixed  $\Delta \in \mathcal{D}$  and  $t_j \rightarrow 0$ .
- C.**  *$\rho$ -Fréchet differentiability.*  $T$  is  $\rho$ -Fréchet differentiable if the differential exists for every  $G_j \in \mathcal{F}_0$  with  $\rho(G_j, F) \rightarrow 0$ .

Since these differentiability conditions are increasingly stringent,  $\rho$ -Fréchet differentiability  $\Rightarrow$   $\rho$ -Hadamard differentiability  $\Rightarrow$  Gâteaux differentiability, but the reverse implications do not hold. Moreover, in the expansion

$$T(G) - T(F) = L_F(G - F) + o(\rho(G, F)) \quad \text{as } \rho(G, F) \rightarrow 0,$$

if the differential  $L_F$  in C exists, it is the same as  $L_F$  in A and B, and if  $L_F$  in B exists, it is the same as in  $L_F$  in A.

We now examine the nature of the Gâteaux differential  $L_F$  of  $T$  by the following heuristics. Due to linearity of  $L_F$

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{T(F + t\Delta) - T(F) - L_F(t\Delta)}{\|t\Delta\|} = 0 \\ & \iff \lim_{t \rightarrow 0} \frac{T(F + t\Delta) - T(F) - tL_F(\Delta)}{t} = 0 \\ & \iff \lim_{t \rightarrow 0} \frac{T(F + t\Delta) - T(F)}{t} = L_F(\Delta), \text{ ie,} \\ & L_F(G - F) = \frac{d}{dt} T((1-t)F + tG)|_{t=0}. \end{aligned}$$

Let  $\delta_x$  be the cdf with its entire mass at  $x$  (ie,  $\delta_x(u) = I_{[x, \infty)}(u)$ ). Then for every  $u$ ,

$$\begin{aligned} & \int \{\delta_x(u) - F(u)\} dG(x) = \int_{x \leq u} dG(x) - F(u) = G(u) - F(u) \text{ and} \\ & \int \{\delta_x(u) - F(u)\} dF(x) = 0, \text{ ie,} \\ & \Delta = G - F = \int (\delta_x - F) d[G(x) - F(x)]. \end{aligned} \tag{2}$$

Thus, if we let

$$T'(x; F) = \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0} = L_F(\delta_x - F),$$

then by linearity of  $L_F$ , using Eq. (2), we should have

$$\begin{aligned} L_F(G - F) &= L_F \int (\delta_x - F) d[G(x) - F(x)] \\ &= \int L_F(\delta_x - F) d[G(x) - F(x)] = \int T'(x; F) d[G(x) - F(x)]. \end{aligned} \quad (3)$$

The first-order differential expansion of  $T(G)$  can now be written as

$$\begin{aligned} T(G) - T(F) &= \int T'(x; F) dG(x) + \text{Rem} \\ &= \int T'(x; F) d[G(x) - F(x)] + \text{Rem}, \end{aligned}$$

where the remainder term  $\rightarrow 0$  as  $G \rightarrow F$ .

Letting  $F_n$  play the role of  $G$ , we now have

$$\begin{aligned} \sqrt{n}[T(F_n) - T(F)] &= \sqrt{n} \int T'(x; F) dF_n(x) + R_n \\ &= n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n. \end{aligned} \quad (4)$$

Since  $\sqrt{n}\|F_n - F\|_\infty = O_P(1)$ , following the analogy of the parametric delta method, *we would expect*  $R_n$  to be  $o_P(1)$  and then the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$  would be the same as that of  $n^{-1/2} \sum_{i=1}^n T'(X_i; F)$ . Interchanging the order of operations  $L_F$  and the integration as in Eq. (3) (which needs justification), *we would expect*

$$\begin{aligned} E_F[T'(X; F)] &= \int T'(x; F) dF(x) = \int L_F(\delta_x - F) dF(x) \\ &= L_F \int (\delta_x - F) dF(x) = L_F(0) = 0, \end{aligned}$$

and if we let

$$\sigma^2(F) = \text{Var}_F[T'(X; F)] = \int \{T'(x; F)\}^2 dF(x),$$

assuming that it exists, then we would have

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + o_P(1) \xrightarrow{\mathcal{L}} N(0, \sigma^2(F)).$$

This is known as the functional delta method.

*Remark 10.2.1.*

1. In the robustness literature [55, 56],  $T'(x; F)$  is called the *influence function* and it is denoted by

$$IF(x; F, T) = T'(x; F) = \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0}.$$

The function  $IF(x; F, T)$  measures the rate at which  $T(F)$  changes when  $F$  is contaminated by  $\delta_x$  with a small probability. The contamination  $\delta_x$  is called *gross-error* and

$$\lambda^* = \sup_x |IF(x; F, T)|$$

is called the *gross-error sensitivity* of  $T$  at  $F$ .

2. Since Gâteaux differentiability is too weak, there is no guarantee that  $R_n = o_P(1)$ , as seen in the following example [57]:

Define  $T(F) = \sum_{x \in [0,1]} [F(x) - F(x-)]^\alpha$ ,  $\alpha > 1$ , as a measure of jumps of  $F$  on  $\mathbb{R}$ . For  $F = U$  (uniform distribution on  $[0, 1]$ ) which has no jumps,

$$\begin{aligned} T'(x; U) &= \frac{d}{dt} T((1-t)U + t\delta_x)|_{t=0} = \frac{d}{dt} \sum_{y \in [0,1]} t^\alpha [\delta_x(y) - \delta_x(y-)]^\alpha|_{t=0} \\ &= \frac{d}{dt} t^\alpha|_{t=0} = \alpha t^{\alpha-1}|_{t=0} = 0 \text{ for } \alpha > 1, \end{aligned}$$

because  $U$  has no jump and  $\delta_x(y)$  has exactly one jump of magnitude 1 at  $y = x$ . Now  $F_n$  has  $n$  jumps of  $1/n$  each with probability 1, so  $T(F_n) = n(1/n)^\alpha = n^{1-\alpha}$ . Hence

$$\sqrt{n}[T(F_n) - T(U)] = \sqrt{n}[n^{1-\alpha} - 0] = n^{3/2-\alpha}$$

with probability 1. Thus the expansion

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n$$

becomes  $n^{3/2-\alpha} = 0 + R_n$  and for  $1 < \alpha < 3/2$ ,  $R_n = o_P(1)$  is false.

We now summarize the “potential” of the Gâteaux differentiability approach based on the above heuristics:

Let  $T'(x; F) = \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0}$ . Then

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n.$$

If  $E_F[T'(X; F)] = 0$ ,  $0 < \text{Var}_F[T'(X; F)] = \sigma^2(F) < \infty$  and if  $R_n = o_P(1)$ , then  $\sqrt{n}[T(F_n) - T(F)] \xrightarrow{P} N(0, \sigma^2(F))$  as  $n \rightarrow \infty$ .

To put the above heuristics to work, the main thing is to demonstrate that  $R_n = o_P(1)$  which can be attempted in one of the two ways.

- I. Use the expansion given in Eq. (4) as a working formula and then carry out the following steps:
  - (i) Calculate  $T'(x; F) = \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0}$  which involves a simple one-variable differentiation.
  - (ii) Check the condition  $E_F[T'(X; F)] = \int T'(x; F) dF(x) = 0$ .
  - (iii) Check the condition  $R_n = o_P(1)$  by examining the particular case.

**(iv)** Calculate  $\sigma^2(F) = \int \{T'(x; F)\}^2 dF(x)$ .

- II. Obtain conditions on  $T$  in terms of  $\rho$ -Hadamard or  $\rho$ -Fréchet differentiability, so that  $R_n = o_P(1)$ .

We first illustrate the first approach by the examples listed in [Section 10.1](#) of this chapter and then go into some theoretical considerations needed to pursue the second approach.

**Example 10.2.1** (The Mean). Here  $T(F) = \int u dF(u)$ ,  $T(F_n) = n^{-1} \sum_{i=1}^n X_i = \bar{X}_n$ . The influence function is

$$\begin{aligned} T'(x; F) &= \frac{d}{dt} \int u d[(1-t)F(u) + t\delta_x(u)]|_{t=0} = x - \int u dF(u) \\ &= x - E_F(X) := x - \xi(F). \end{aligned}$$

In this case, the first-order approximation is an identity, because the function  $T(F)$  is already linear. If the sample space is  $\mathbb{R}$ , then the influence function is unbounded, so  $\bar{X}_n$  is not robust.

**Example 10.2.2** (The  $k$ th Central Moment). Let  $\mu_k = T(F) = \int [u - \xi(F)]^k dF(u)$ , where  $\xi(F) = E_F(X)$ . The influence function is

$$\begin{aligned} T'(x; F) &= \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0} \\ &= \frac{d}{dt} \int [u - \xi(F) - t(x - \xi(F))]^k d[F(u) + t(\delta_x(u) - F(u))]|_{t=0} \\ &= -k(x - \xi(F))\mu_{k-1} + (x - \xi(F))^k - \mu_k, \end{aligned}$$

which is unbounded. By routine calculations, we have

$$\begin{aligned} E_F[T'(X; F)] &= 0 \text{ and} \\ \sigma^2(F) &= \mu_{2k} - \mu_k^2 - 2k\mu_{k-1}\mu_{k+1} + k^2\mu_{k-1}\mu_2. \end{aligned}$$

Finally, the remainder term of this one-term expansion of  $\sqrt{n}[T(F_n) - T(F)]$  is

$$\begin{aligned} R_n &= \sqrt{n}[T(F_n) - T(F)] - n^{-1/2} \sum_{i=1}^n T'(X_i; F) \\ &= \sqrt{n} \left[ n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^k - \mu_k \right] \\ &\quad - n^{-1/2} \sum_{i=1}^n \left[ (X_i - \xi(F))^k - \mu_k - k\mu_{k-1}(X_i - \xi(F)) \right] \\ &= n^{-1/2} \sum_{i=1}^n \left[ (X_i - \bar{X}_n)^k - (X_i - \xi(F))^k \right] + k\mu_{k-1}\sqrt{n}(\bar{X}_n - \xi(F)) \\ &= n^{-1/2} \sum_{i=1}^n \left[ (Y_i - \bar{Y}_n)^k - Y_i^k \right] + k\mu_{k-1}\sqrt{n}\bar{Y}_n \end{aligned} \tag{5}$$

with  $Y_i = X_i - \xi(F)$ . Then  $E_F[Y_i] = 0$  and  $E_F[Y_i^{k-j}] = \mu_{k-j}$ . Now expand the last expression in Eq. (5), make some algebraic rearrangements and note that  $\sqrt{n}\bar{Y}_n = O_P(1)$ ,  $n^{-1} \sum_{i=1}^n Y_i^{k-j} = \mu_{k-j} + o_P(1)$  for  $j = 2, \dots, k$ . This shows that  $R_n = o_P(1)$ .

**Example 10.2.3** (The  $p$ -Quantile). For  $0 < p < 1$ , the  $p$ -quantile of  $F$  is

$$T(F) = F^{-1}(p) = \inf\{x: F(x) \geq p\}, \quad \text{so that } F(F^{-1}(p)) = p.$$

Let  $F_t = (1-t)F + t\delta_x$  for a give  $x$ . Now differentiating both sides of the identity

$$p = F_t(F_t^{-1}(p)) = (1-t)F(F_t^{-1}(p)) + t\delta_x(F_t^{-1}(p))$$

with respect to  $t$ , evaluated at 0, we have.

$$\begin{aligned} 0 &= \left[ \left\{ -F(F_t^{-1}(p)) + (1-t)f(F_t^{-1}(p)) \frac{d}{dt} F_t^{-1}(p) \right\} + \left\{ \delta_x(F_t^{-1}(p)) + t \frac{d}{dt} \delta_x(F_t^{-1}(p)) \right\} \right]_{t=0} \\ &= -F(F^{-1}(p)) + f(F^{-1}(p)) \frac{d}{dt} F_t^{-1}(p) \Big|_{t=0} + I_{[x, \infty)}(F^{-1}(p)) + 0 \\ &= f(F^{-1}(p)) \frac{d}{dt} F_t^{-1}(p) \Big|_{t=0} - \{p - I_{(-\infty, F^{-1}(p))}(x)\}. \end{aligned}$$

Hence

$$T'(x; F) = \frac{d}{dt} F_t^{-1}(p) \Big|_{t=0} = \frac{p - I_{(-\infty, F^{-1}(p))}(x)}{f(F^{-1}(p))}, \quad (6)$$

and the leading term of the expression of  $\sqrt{n}[T(F_n) - T(F)]$  is

$$\begin{aligned} n^{-1/2} \sum_{i=1}^n T'(X_i; F) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{p - I_{(-\infty, F^{-1}(p))}(X_i)}{f(F^{-1}(p))} \\ &\xrightarrow{\mathcal{L}} N\left(0, \frac{p(1-p)}{f^2(F^{-1}(p))}\right), \end{aligned}$$

provided that the remainder term  $R_n = o_P(1)$ . Indeed,  $R_n = O_{a.s.}(n^{-1/4} \log n)$ , ie, there exists  $C$  such that  $P[|R_n| > Cn^{-1/4} \log n \text{ i.o.}] = 0$ . For a proof, Bahadur [58].

*Remark 10.2.2.*

1. The expression (6), rewritten as

$$X_{n:[np]} = F^{-1}(p) + \frac{1}{n} \sum_{i=1}^n \frac{p - I_{(-\infty, F^{-1}(p))}(X_i)}{f(F^{-1}(p))} + R_n$$

is known as Bahadur representation of a sample quantile.

2. Since  $\left[ p - I_{(-\infty, F^{-1}(p)]}(X_i) \right] / f(F^{-1}(p))$ ,  $i = 1, \dots, n$  are iid with mean zero and variance  $p(1-p)/f^2(F^{-1}(p))$ ,  $X_{n:[np]} \xrightarrow{a.s.} F^{-1}(p)$ .
3. The Bahadur representation can be obtained by a simple heuristic argument. Let  $F_n$  be the edf of  $X_1, \dots, X_n$ . Then  $F_n(X_{n:[np]}) = [np]/n$ . Now note that if  $\xi \leq X_{n:[np]}$ , then

$$\begin{aligned} (X_{n:[np]} - \xi)f(\xi) &\approx F(X_{n:[np]}) - F(\xi) \approx F_n(X_{n:[np]}) - F_n(\xi) \\ &= [np] - n^{-1} \sum_{i=1}^n I_{(-\infty, \xi]}(X_i) \approx n^{-1} \sum_{i=1}^n [p - I_{(-\infty, \xi]}(X_i)] \end{aligned}$$

and similarly, for  $X_{n:[np]} \leq \xi$ ,

$$(\xi - X_{n:[np]})f(\xi) \approx n^{-1} \sum_{i=1}^n [I_{(-\infty, \xi]}(X_i) - p].$$

In both cases,

$$X_{n:[np]} \approx \xi + n^{-1} \sum_{i=1}^n [p - I_{(-\infty, \xi]}(X_i)]/f(\xi),$$

where  $\xi = F^{-1}(p)$ .

4. Asymptotic joint distributions of several sample quantiles can be obtained similarly.

**Example 10.2.4** (The  $V$ -Statistic). Let  $g: \mathbb{R}^r \rightarrow \mathbb{R}$  (or more generally,  $g: \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}$ ) be such that  $g(y_{i_1}, \dots, y_{i_r}) = g(y_1, \dots, y_r)$  for all  $(y_1, \dots, y_r)$  and all permutations  $(i_1, \dots, i_r)$  of  $(1, \dots, r)$ . Then  $T(F) = \int \dots \int g(y_1, \dots, y_r) \prod_{i=1}^r dF(y_i)$  for  $F \in \mathcal{F}_0 = \{F: E_F[|g(Y_1, \dots, Y_r)|] < \infty\}$  is called a  $V$ -functional and its corresponding  $V$ -statistic is

$$T(F_n) = n^{-r} \sum_{i_1=1}^n \dots \sum_{i_r=1}^n g(X_{i_1}, \dots, X_{i_r})$$

based on a random sample  $X_1, \dots, X_n$  from  $F$ . Since  $g$  is symmetric in its coordinates,

$$\begin{aligned} T((1-t)F + t\delta_x) &= \int \dots \int g(y_1, \dots, y_r) \prod_{i=1}^r d[(1-t)F(y_i) + t\delta_x(y_i)] \\ &= (1-t)^r T(F) + \sum_{c=1}^r t^c (1-t)^{r-c} r^{(c)} \\ &\quad \times \int \dots \int g(x, \dots, x, y_{c+1}, \dots, y_r) \prod_{i=c+1}^r dF(y_i) \end{aligned}$$

which is a polynomial in  $t$ , so its derivative with respect to  $t$ , evaluated at 0 is the coefficient of  $t$ . Thus

$$\begin{aligned}
T'(x; F) &= \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0} \\
&= r \int \cdots \int g(x, y_2, \dots, y_r) \prod_{i=2}^r dF(y_i) - rT(F) \\
&= r\{\text{E}_F[g(X_1, X_2, \dots, X_r)|X_1 = x] - T(F)\} = rh_1(x), \text{ where} \\
h_1(x) &= \text{E}_F[g(X_1, X_2, \dots, X_r)|X_1 = x] - T(F).
\end{aligned}$$

Clearly,

$$r^{-1}\text{E}_F[T'(X; F)] = 0 \text{ and } \sigma^2(F) = r^2\text{E}_F[h_1^2(X)].$$

Hence subject to the verification of the remainder term

$$R_n = \sqrt{n}[T(F_n) - T(F)] - rn^{-1/2} \sum_{i=1}^n h_1(X_i)$$

being  $o_P(1)$ , we have

$$\sqrt{n}[T(F_n) - T(F)] \xrightarrow{\mathcal{L}} N(0, \sigma^2(F)).$$

Examine the discrepancy between the  $V$ -statistic and the corresponding  $U$ -statistic to see that this result is what one would expect.

### 10.3 The $L$ -Estimators

Let  $F(\cdot - \theta)$  be a cdf with pdf  $f(\cdot - \theta)$  where  $f$  is symmetric about 0 and  $\theta \in \mathbb{R}$  is unknown. Then  $\theta$  is a location parameter which is the point of symmetry of the unknown distribution, which is the median of  $F(\cdot - \theta)$  and also the mean if it exists. If  $X_1, \dots, X_n$  is a random sample from  $F(\cdot - \theta)$ , then the sample mean  $\bar{X}_n$ , being the UMVUE and the MLE of  $\theta$  if  $F(\cdot - \theta)$  is normal with mean  $\theta$ , is a very good estimator, but it is not so good if  $F(\cdot - \theta)$  is Cauchy with median  $\theta$ . The median  $N_{n:[n/2]}$  is no good in case of normal distributions, but does not break down like  $\bar{X}_n$  for Cauchy distributions.

The reason for  $\bar{X}_n$  performing so poorly for the Cauchy distribution is due to its heavy tails. The Cauchy pdf tends to 0 at the rate of  $1/x^2$  as  $x \rightarrow \pm\infty$  as opposed to the  $e^{-x^2/2}$  rate for the normal pdf. This makes the extreme-order statistics unstable. The sample mean can also be viewed as  $n^{-1} \sum_{i=1}^n X_{n:i}$ , where  $X_{n:1} < \dots < X_{n:n}$  are the order statistics in  $(X_1, \dots, X_n)$ . To protect  $\bar{X}_n$  from being drastically affected by possibly heavy tails of the underlying distribution, it would seem reasonable to redistribute the weights on the order statistics so that the extreme ones are de-emphasized. This leads to the consideration of estimators which are linear functions of order statistics. These estimators include the sample mean, the sample median, and also a class of estimators called  $\alpha$ -trimmed means with  $0 < \alpha < 1/2$  defined by  $\bar{X}_{n:(\alpha)} = (n - 2[\alpha n])^{-1} \sum_{i=[\alpha n]+1}^{n-\lceil \alpha n \rceil} X_{n:i}$ , which fall between the sample mean and the sample median.

**Definition 10.3.1.** An  $L$ -functional  $T: \mathcal{F}_0 \rightarrow \mathbb{R}$  is defined as

$$T(F) = \int_{-\infty}^{\infty} xJ[F(x)] dF(x) = \int_0^1 F^{-1}(u)J(u) du,$$

$F \in \mathcal{F}_0$  and  $J: [0, 1] \rightarrow R$ ,

and if  $X_1, \dots, X_n$  is a random sample from  $F \in \mathcal{F}_0$  with edf  $F_n$ , then

$$T(F_n) = \int_{-\infty}^{\infty} xJ[F_n(x)] dF_n(x) = n^{-1} \sum_{i=1}^n J(i/n)X_{n:i}$$

is called the  $L$ -estimator of  $T(F)$ . The function  $J$  is called the score function.

We now calculate the influence function and the remainder  $R(G, F)$  in the expansion

$$T(G) - T(F) = \int T'(x; F) d[G(x) - F(x)] + R(G, F).$$

Differentiating

$$T(F + t(\delta_x - F)) = \int uJ[F(u) + t(\delta_x(u) - F(u))]d[F(u) + t(\delta_x(u) - F(u))]$$

with respect to  $t$  at  $t = 0$ , we have

$$T'(x; F) = \int u(\delta_x(u) - F(u))J'[F(u)] dF(u) + \int uJ[F(u)] d(\delta_x(u) - F(u)).$$

Integrate the second integral by parts and make some algebraic rearrangements to get

$$\begin{aligned} T'(x; F) &= - \int (\delta_x(u) - F(u))J[F(u)] du \\ &= \int_{-\infty}^x J[F(u)] du - \int_{-\infty}^{\infty} (1 - F(u))J[F(u)] du. \end{aligned}$$

The remainder term is

$$R(G, F) = T(G) - T(F) - \int T'(x; F) d[G(x) - F(x)],$$

where

$$\begin{aligned} T(G) - T(F) &= \int_0^1 [G^{-1}(u) - F^{-1}(u)]J(u) du = \int_0^1 \left[ \int_{F^{-1}(u)}^{G^{-1}(u)} dx \right] J(u) du \\ &= \int_0^1 \int_{-\infty}^{\infty} I_{[F^{-1}(u), G^{-1}(u)]}(x)J(u) dx du \\ &= \int_{-\infty}^{\infty} \int_0^1 I_{[G(x), F(x)]}(u)J(u) du dx \\ &= \int_{-\infty}^{\infty} \left[ \int_{G(x)}^{F(x)} J(u) du \right] dx = - \int_{-\infty}^{\infty} \left[ \int_{F(x)}^{G(x)} J(u) du \right] dx, \text{ and} \\ \int T'(x; F) d[G(x) - F(x)] &= \int_{-\infty}^{\infty} \left[ \int_{-\infty}^x J[F(u)] du \right] d[G(x) - F(x)] - 0 \\ &= - \int_{-\infty}^{\infty} [G(x) - F(x)]J[F(x)] dx. \end{aligned}$$

Thus

$$\begin{aligned} R(G, F) &= - \int_{-\infty}^{\infty} \left[ \int_{F(x)}^{G(x)} J(u) du - (G(x) - F(x))J(F(x)) \right] dx \\ &= - \int_{-\infty}^{\infty} W_G(x)[G(x) - F(x)] dx, \text{ where} \\ W_G(x) &= (G(x) - F(x))^{-1} \int_{F(x)}^{G(x)} J(u) du - J(F(x)). \end{aligned} \quad (7)$$

Now letting  $F_n$  play the role of  $G$ , we have

$$\begin{aligned} \sqrt{n}[T(F_n) - T(F)] &= n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n, \text{ where} \\ T'(x; F) &= - \int (\delta_x(u) - F(u))J[F(u)] du \text{ and } R_n = R(F_n, F). \end{aligned} \quad (8)$$

We now verify that  $E_F[T'(X; F)] = 0$  and obtain a formula for  $\sigma^2(F) = \text{Var}_F[T'(X; F)]$ . This will then give us the asymptotic distribution of the  $L$ -estimator, namely,

$$\sqrt{n}[T(F_n) - T(F)] \xrightarrow{\mathcal{L}} N(0, \sigma^2(F))$$

subject to  $R_n = R(F_n, F) = o_P(1)$ .

First,

$$\begin{aligned} E_F[T'(X; F)] &= - \int \left[ \int (\delta_x(u) - F(u))J[F(u)] du \right] dF(x) \\ &= - \int \left[ \int (\delta_x(u) - F(u)) dF(x) \right] J(F(u)) du = 0, \end{aligned}$$

because  $\int (\delta_x(u) - F(u)) dF(x) = 0$  for all  $u$ .

Next,

$$\begin{aligned} \sigma^2(F) &= \text{Var}_F[T'(X; F)] \\ &= \int \left[ \int (\delta_x(u) - F(u))J[F(u)] du \right]^2 dF(x) \\ &= 2 \iint_{u < v} \left[ \int (\delta_x(u) - F(u))(\delta_x(v) - F(v)) dF(x) \right] J(F(u))J(F(v)) du dv \\ &= 2 \iint_{u < v} F(u)(1 - F(v))J(F(u))J(F(v)) du dv \\ &= \iint [F(\min(u, v)) - F(u)F(v)]J(F(u))J(F(v)) du dv, \end{aligned} \quad (9)$$

because

$$\begin{aligned} \int (\delta_x(u) - F(u))(\delta_x(v) - F(v)) dF(x) &= (1 - F(u))(1 - F(v))F(u) \\ &\quad - F(u)(1 - F(v))(F(v) - F(u)) + F(u)F(v)(1 - F(v)) \\ &= F(u)(1 - F(v)). \end{aligned}$$

An alternate formula for  $\sigma^2(F)$  can be obtained by first replacing  $(u, v)$  by  $(x, y)$  and then letting  $F(x) = u$  and  $F(y) = v$  in Eq. (9). Thus

$$\begin{aligned}\sigma^2(F) &= 2 \iint_{x < y} F(x)(1 - F(y))J(F(x))J(F(y)) dx dy \\ &= 2 \iint_{u < v} u(1 - v) \left\{ J(u)/f(F^{-1}(u)) \right\} \left\{ J(v)/f(F^{-1}(v)) \right\} du dv \\ &= 2 \iint_{u < v} u(1 - v)A'(u)A'(v) du dv, \text{ where} \\ A(u) &= \int_0^u \frac{J(t)}{f(F^{-1}(t))} dt + c\end{aligned}$$

where  $c$  is such that  $A(1) = 0$ , making  $A'(u) = J(u)/f(F^{-1}(u))$ . Thus

$$\begin{aligned}\sigma^2(F) &= 2 \int_{u=0}^1 \left[ \int_{v=u}^1 A'(v) dv \right] uA'(u) du - \int_0^1 \int_0^1 uvA'(u)A'(v) du dv \\ &= 2 \int_0^1 [0 - A(u)]uA'(u) du - \left( \int_0^1 uA'(u) du \right)^2 \\ &= - \int_0^1 u dA^2(u) - \left( \int_0^1 u dA(u) \right)^2 = \int_0^1 A^2(u) du - \left( \int_0^1 A(u) du \right)^2,\end{aligned}\quad (10)$$

using integration by parts in both integrals.

### 10.3.1 Asymptotic Distribution of $\alpha$ -Trimmed Mean When $f$ Is Symmetric About $\theta$

The  $L$ -functional for  $\alpha$ -trimmed mean is:  $T(F) = (1 - 2\alpha)^{-1} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x)$  (ie, the score function is  $J(t) = (1 - 2\alpha)^{-1} I_{[\alpha, 1-\alpha]}(t)$ ). Since  $f = F'$  is symmetric about  $\theta$ ,  $F^{-1}(1 - \alpha) - \theta = \theta - F^{-1}(\alpha)$ , so that

$$\int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} (x - \theta) dF(x) = 0 \text{ and } T(F) = (1 - 2\alpha)^{-1} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x) = \theta.$$

Next,

$$\begin{aligned}A(u) &= \int_0^u \frac{(1 - 2\alpha)^{-1} I_{[\alpha, 1-\alpha]}(t)}{f(F^{-1}(t))} dt \\ &= \begin{cases} 0 & \text{if } 0 < u < \alpha \\ (1 - 2\alpha)^{-1} [F^{-1}(u) - F^{-1}(\alpha)] & \text{if } \alpha \leq u \leq 1 - \alpha \\ (1 - 2\alpha)^{-1} [F^{-1}(1 - \alpha) - F^{-1}(\alpha)] & \text{if } 1 - \alpha < u < 1. \end{cases}\end{aligned}$$

To calculate  $\int_0^1 A(t) dt$  and  $\int_0^1 A^2(t) dt$ , use that due to symmetry of  $f = F'$  about  $\theta$ ,  $F^{-1}(1 - \alpha) - \theta = \theta - F^{-1}(\alpha)$ . Using these facts, we obtain, after algebraic simplifications,

$$\begin{aligned}\sigma^2(F) &= \int_0^1 A^2(u) du - \left( \int_0^1 A(u) du \right)^2 \\ &= \frac{1}{(1-2\alpha)^2} \left[ 2\alpha(\theta - F^{-1}(\alpha))^2 + \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} (x-\theta)^2 dF(x) \right],\end{aligned}$$

by Eq. (10).

## 10.4 The $M$ -Estimators

The method of maximum likelihood is based on the fact that subject to identifiability of the family  $\{f(\cdot, t)\}$ , if  $X$  has pdf  $f(x, \theta)$ , then the function  $E_\theta[\log f(X, t)]$  has  $t = \theta$  as its unique maximizer. The MLE of  $\theta$  based on a random sample  $X_1, \dots, X_n$  is the maximizer of  $n^{-1} \sum_{i=1}^n \log f(X_i, t)$  which is a natural estimate of  $E_\theta[\log f(X, t)]$ . Equivalently, the MLE of  $\theta$  is a solution of  $\sum_{i=1}^n \Psi(X_i, t) = 0$ , where  $\Psi(x, t) = \partial \log f(x, t) / \partial t$ , assuming that  $f$  is smooth. In the location problem,

$$\Psi(x, t) = \partial \log f(x-t) / \partial t = -\frac{f'}{f}(x-t),$$

so the MLE of a location parameter is the solution of  $\sum_{i=1}^n \Psi(X_i - t) = 0$  where  $\Psi = -f'/f$ .

The MLEs have good properties under the correct model, but if not, then  $\Psi(x, t)$  may be very unstable for some  $x$ , as in the location problem in which  $-(f'/f)(x) = x$  for  $N(0, 1)$  and  $X_i - t$  for extreme observations are unstable if the true distribution is heavy-tailed.

The  $M$ -estimators attempt to overcome this weakness of MLEs by using a function  $\rho(x-t)$  in the location problem instead of, but somewhat similar to  $\log f(x-t)$ , and then solving for  $t$  in the equation  $\sum_{i=1}^n \Psi(X_i - t) = 0$ , where  $\Psi = \rho'$  instead of  $-f'/f$ . In general, we replace  $-\log f(x, t)$  by  $\rho(x, t)$  and solve for  $t$  in the equation  $\sum_{i=1}^n \Psi(X_i, t) = 0$ , where  $\Psi(x, t) = \partial \rho(x, t) / \partial t$ .

We now formally define an  $M$ -functional  $T(F)$ , of which  $T(F_n)$  will be an  $M$ -estimator.

**Definition 10.4.1.** Let  $\rho: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$  and let  $\Theta$  be an open subset of  $\mathbb{R}$ . Then

$$T(F) = \arg \min_{t \in \Theta} \int \rho(x, t) dF(x), \quad F \in \mathcal{F}_0$$

is an  $M$ -functional, and if  $F_n$  is the edf of a random sample  $X_1, \dots, X_n$  from  $F$ , then

$$T(F_n) = \arg \min_{t \in \Theta} \sum_{i=1}^n \rho(X_i, t)$$

is the  $M$ -estimator of  $T(F)$ .

If  $\Psi(x, t) = \partial \rho(x, t) / \partial t$  exists and if

$$\lambda_F(t) = \int \Psi(x, t) dF(x) = \frac{\partial}{\partial t} \int \rho(x, t) dF(x), \quad (11)$$

then  $T(F)$  is a solution of  $\lambda_F(t) = 0$ , and so

$$\lambda_F(T(F)) = \int \Psi(x, T(F)) dF(x) = 0 \quad \text{for all } F \in \mathcal{F}_0. \quad (12)$$

The  $M$ -estimator  $T(F_n)$  of  $T(F)$  can be equivalently be expressed as a solution of the equation  $\sum_{i=1}^n \Psi(X_i, t) = 0$ .

To calculate  $T'(x; F) = \frac{d}{dt} T((1-t)F + t\delta_x)|_{t=0}$ , start with

$$\begin{aligned} 0 &= \lambda_{(1-t)F+t\delta_x}((1-t)F + t\delta_x) \\ &= \int \Psi(u, T((1-t)F + t\delta_x)) d[(1-t)F(u) + t\delta_x(u)] \\ &= (1-t) \int \Psi(u, T((1-t)F + t\delta_x)) dF(u) + t\Psi(x, T((1-t)F + t\delta_x)), \end{aligned}$$

and differentiate both sides with respect to  $t$  at  $t = 0$  to obtain

$$\begin{aligned} 0 &= - \int \Psi(u, T(F)) dF(u) + \frac{d}{dt} \int \Psi(u, T((1-t)F + t\delta_x)) dF(u)|_{t=0} \\ &\quad + \Psi(x, T(F)) + 0 \\ &= -\lambda_F(T(F)) + \frac{d}{dt} \lambda_F(T((1-t)F + t\delta_x))|_{t=0} + \Psi(T(F)) \\ &= 0 + \lambda'_F(T(F))T'(x; F) + \Psi(x, T(F)), \end{aligned}$$

using Eqs. (11) and (12).

Assuming  $\lambda'_F(T(F)) \neq 0$ , we thus have

$$T'(x; F) = -\Psi(x, T(F))/\lambda'_F(T(F)). \quad (13)$$

We now check that

$$E_F[T'(X; F)] = -E_F[\Psi(X, T(F))]/\lambda'_F(T(F)) = -\lambda_F(T(F))/\lambda'_F(T(F)) = 0$$

by Eqs. (11) and (12), and

$$\sigma^2(F) = \sigma^2(F, \Psi) = \text{Var}_F[T'(X; F)] = \frac{\int \Psi^2(x, T(F)) dF(x)}{\{\lambda'_F(T(F))\}^2}. \quad (14)$$

We can now conclude that

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n \xrightarrow{\mathcal{L}} N(0, \sigma^2(F, \Psi)),$$

subject to the verification  $R_n = o_P(1)$ .

The MLE is also an  $M$ -estimator with  $\Psi(x, t) = -\partial \log f(x, t)/\partial t$ , and letting  $F_\theta = F(\cdot, \theta)$ ,  $T(F_\theta) = \theta$ . Now

$$\begin{aligned} \lambda'_{F_\theta}(T(F_\theta)) &= \lambda'_{F_\theta}(\theta) = \int \left( -\frac{\partial^2 \log f(x, t)}{\partial t^2} \right)_{t=\theta} f(x, \theta) dx = I_f(\theta), \text{ and} \\ \int \Psi^2(x, T(F_\theta)) dF_\theta(x) &= \int \left( \frac{\partial \log f(x, t)}{\partial t} \right)_{t=\theta}^2 f(x, \theta) dx = I_f(\theta), \end{aligned}$$

where  $I_f(\theta)$  is the Fisher-information of the family  $\{f(x, t), t \in \Theta\}$  at  $t = \theta$ . Thus  $\sigma^2(F_\theta, \Psi) = 1/I_f(\theta)$  by Eq. (14).

In the location problem,  $f(x, \theta) = f(x - \theta)$  where  $f$  is a symmetric (about 0) pdf. Here

$$\frac{\partial \log f(x-t)}{\partial t} \Big|_{t=\theta} = -\frac{f'(x-\theta)}{f(x-\theta)} \text{ and}$$

$$I_f(\theta) = \int \left\{ \frac{f'(x-\theta)}{f(x-\theta)} \right\}^2 f(x-\theta) dx = \int \left\{ \frac{f'(x)}{f(x)} \right\}^2 f(x) dx \text{ for all } \theta.$$

#### 10.4.1 A Minimax Approach to the Choice of $\Psi$

For robust estimation of a location parameter  $\theta$ , instead of assuming  $F$  to be a known cdf, we work within the model

$$\mathcal{F} = \mathcal{F}(G, \varepsilon) = \left\{ F = (1 - \varepsilon')G + \varepsilon'H : H \text{ is a cdf symmetric about 0}, \varepsilon' \in [0, \varepsilon] \right\},$$

where  $G$  is a specified cdf which is symmetric about 0 and  $\varepsilon$  is a specified positive number (ie, we assume  $F$  to be a symmetric cdf lying “within  $\varepsilon$  distance” of a specified symmetric cdf  $G$ ). For given  $G$  and  $\varepsilon$ , we now look for  $\Psi_0 = \Psi_0(G, \varepsilon)$  such that

$$\sup_{F \in \mathcal{F}(G, \varepsilon)} \sigma^2(F, \Psi_0(G, \varepsilon)) \leq \sup_{F \in \mathcal{F}(G, \varepsilon)} \sigma^2(F, \Psi) \quad \text{for all } \Psi.$$

Then the  $M$ -estimator with  $\Psi = \Psi_0(G, \varepsilon)$  is minimax for the family  $\mathcal{F}(G, \varepsilon)$  in the sense of minimizing the maximum possible asymptotic variance.

For  $G = \Phi$ , the cdf of  $N(0, 1)$ , the solution of this minimax problem lies in the class of “Huber Functions”:

$$\Psi_0(x) = \begin{cases} -k & x \leq -k \\ x & |x| < k \\ k & x \geq k, \end{cases}$$

where  $k$  is given in terms of  $\varepsilon$  by the formula

$$\int_{-k}^k \phi(x) dx + (2/k)\phi(k) = 1/(1 - \varepsilon), \quad \phi = \Phi' = \text{pdf of } N(0, 1).$$

See Huber [56].

#### An Alternative Derivation of the Asymptotic Distribution of $M$ -Estimators When the Score Function $\Psi$ Is Monotone

Let  $\hat{\theta}_n$  be an  $M$ -estimator of the location parameter  $\theta$  of a family of pdf's  $\{f(x, \theta) = f(x - \theta), \theta \in \mathbb{R}\}$  where  $f$  is symmetric about 0 and suppose that the score function is antisymmetric and monotone increasing.

Since  $\hat{\theta}_n$  is the solution of  $\sum_{i=1}^n \Psi(X_i - t) = 0$  and  $\sum_{i=1}^n \Psi(X_i - t)$  is monotone decreasing in  $t$ ,  $\sum_{i=1}^n \Psi(X_i - t) \leq 0$  for all  $t \geq \hat{\theta}_n$  and  $\sum_{i=1}^n \Psi(X_i - t) \geq 0$  for all  $t \leq \hat{\theta}_n$ . Hence

$$\begin{aligned}
P_\theta \left[ \sqrt{n}(\hat{\theta}_n - \theta) \leq a \right] &= P_\theta \left[ \theta + a/\sqrt{n} \geq \hat{\theta}_n \right] \\
&= P_\theta \left[ n^{-1/2} \sum_{i=1}^n \Psi(X_i - \theta - a/\sqrt{n}) \leq 0 \right] \\
&= P_0 \left[ n^{-1/2} \sum_{i=1}^n \Psi(X_i - a/\sqrt{n}) \leq 0 \right].
\end{aligned}$$

Now write

$$n^{-1/2} \sum_{i=1}^n \Psi(X_i - a/\sqrt{n}) = n^{-1/2} \sum_{i=1}^n \Psi(X_i) - an^{-1} \sum_{i=1}^n \Psi'(X_i) + R_n.$$

Since  $E_0[\Psi(X_i)] = 0$  due to  $\Psi$  being antisymmetric,

$$n^{-1/2} \sum_{i=1}^n \Psi(X_i) \xrightarrow{\mathcal{L}} N(0, \text{Var}_0[\Psi(X_i)])$$

provided that  $\text{Var}_0[\Psi(X_i)] = \int \Psi^2(x)f(x) dx < \infty$ . Also,

$$n^{-1} \sum_{i=1}^n \Psi'(X_i) = E_0[\Psi'(X_i)] + o_P(1)$$

provided that  $E_0[\Psi'(X_i)] = \int \Psi'(x)f(x) dx$  exists.

Thus if  $R_n = o_P(1)$  by regularity conditions on  $F$  and  $\Psi$ , then

$$n^{-1/2} \sum_{i=1}^n \Psi(X_i - a/\sqrt{n}) \xrightarrow{\mathcal{L}} N(-aE_0[\Psi'(X)], \text{Var}_0[\Psi(X)]).$$

Consequently,

$$\begin{aligned}
\lim_{n \rightarrow \infty} P_\theta \left[ \sqrt{n}(\hat{\theta}_n - \theta) \leq a \right] &= \lim_{n \rightarrow \infty} P_0 \left[ n^{-1/2} \sum_{i=1}^n \Psi(X_i - a/\sqrt{n}) \leq 0 \right] \\
&= \Phi \left( \frac{aE_0[\Psi'(X)]}{\sqrt{\text{Var}_0[\Psi'(X)]}} \right) = \Phi \left( \frac{a}{\sqrt{\text{Var}_0[\Psi(X)]/E_0^2[\Psi'(X)]}} \right),
\end{aligned}$$

provided that  $E_0[\Psi'(X)] \neq 0$ , that is,

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} N(0, \sigma^2(F, \Psi)), \text{ where}$$

$$\sigma^2(F, \Psi) = \frac{\text{Var}_0[\Psi(X)]}{E_0^2[\Psi'(X)]} = \frac{\int \Psi^2(x)f(x) dx}{[\int \Psi'(x)f(x) dx]^2}.$$

In the above,  $\Psi$  is assumed to be strictly increasing.

If  $\Psi$  is merely nondecreasing, then the equation  $\sum_{i=1}^n \Psi(X_i - t) = 0$  may be satisfied on an entire interval. In that case, the above argument would need minor modification.

## 10.5 A Relation Between $L$ -Estimators and $M$ -Estimators

Let  $F$  be a cdf, symmetric about 0 and let  $f = F'$  be the pdf. For simplicity, assume  $F$  to be strictly increasing.

Consider an  $M$ -estimator of  $\theta$  based on iid observations  $X_1, \dots, X_n$  from  $F(x - \theta)$ , with score function  $\Psi$  having  $\Psi(-x) = -\Psi(x)$  for all  $x$ . Note that nothing changes if  $\Psi$  is multiplied by a constant, so we normalize  $\Psi$  to make  $\int \Psi'(x)f(x) dx = 1$  (assuming that  $\int \Psi'(x)f(x) dx \neq 0$ ). Then the asymptotic variance of the  $M$ -estimator with score function  $\Psi$  is

$$\sigma^2(F, \Psi) = \frac{\int \Psi^2(x)f(x) dx}{[\int \Psi'(x)f(x) dx]^2} = \int \Psi^2(x)f(x) dx.$$

Now consider an  $L$ -estimator with score function

$$J_\Psi(u) = \Psi'(F^{-1}(u)), \quad 0 < u < 1.$$

Then the asymptotic variance of this  $L$ -estimator is the same as  $\sigma^2(F, \Psi)$ . To see this, note that here

$$A'_\Psi(u) = J_\Psi(u)/f(F^{-1}(u)) = \Psi'(F^{-1}(u))/f(F^{-1}(u)),$$

so that

$$\begin{aligned} A_\Psi(t) &= \int_{1/2}^t A'_\Psi(u) du = \int_{1/2}^t \left\{ \Psi'(F^{-1}(u))/f(F^{-1}(u)) \right\} du \\ &= \int_{F^{-1}(1/2)}^{F^{-1}(t)} \left\{ \Psi'(x)/f(x) \right\} dF(x) = \int_0^{F^{-1}(t)} \Psi'(x) dx \\ &= \Psi(F^{-1}(t)) - \Psi(0) = \Psi(F^{-1}(t)). \end{aligned}$$

Hence

$$\int_0^1 A_\Psi(t) dt = \int_0^1 \Psi(F^{-1}(t)) dt = \int_{-\infty}^{\infty} \Psi(x)f(x) dx = 0 \text{ and}$$

$$\int_0^1 A_\Psi^2(t) dt = \int_0^1 \left\{ \Psi(F^{-1}(t)) \right\}^2 dt = \int_{-\infty}^{\infty} \Psi^2(x)f(x) dx.$$

Thus the asymptotic variance of the  $L$ -estimator with score function  $J_\Psi$  is

$$\int_0^1 A_\Psi^2(t) dt - \left( \int_0^1 A_\Psi(t) dt \right)^2 = \int_{-\infty}^{\infty} \Psi^2(x)f(x) dx = \sigma^2(F, \Psi).$$

## 10.6 The Remainder Term $R_n$

Throughout, we consider the case of  $d = 1$ , that is,  $X_1, \dots, X_n$  are real-valued iid rv's with cdf  $F$  and  $F_n$  is the edf of  $X_1, \dots, X_n$ .

We shall consider the following distances in

$\mathcal{F}$  = collection of all cdf's on  $\mathbb{R}$ , and

$\mathcal{F}_1$  = collection of all cdf's on  $\mathbb{R}$  with finite mean:

Sup-norm distance:  $\rho_\infty(F_1, F_2) = \sup_x |F_1(x) - F_2(x)|$  for all  $F_1$  and  $F_2 \in \mathcal{F}$ .

$L_p$ -distance:  $\rho_{L_p}(F_1, F_2) = [\int_{\mathbb{R}} |F_1(x) - F_2(x)|^p dx]^{1/p}$ ,  $p \geq 1$  for all  $F_1$  and  $F_2 \in \mathcal{F}_1$ .

**Lemma 10.6.1.** (i)  $\rho_\infty(F_n, F) = O_P(n^{-1/2})$ , (ii)  $\rho_{L_p}(F_n, F) = O_P(n^{-1/2})$  if either  $1 \leq p < 2$  and  $\int [F(x)(1 - F(x))]^{p/2} dx < \infty$  or  $p \geq 2$ .

*Proof.* By the DKW Theorem (See [Theorem 9.2.4](#)),

$$\begin{aligned} E[\sqrt{n}\rho_\infty(F_n, F)] &= \int_0^\infty P[\sqrt{n}\rho_\infty(F_n, F) > y] dy \\ &\leq \int_0^\infty C \exp[-2n(y/\sqrt{n})^2] dy \\ &= C \int_0^\infty \exp[-2y^2] dy \quad \text{for all } n. \end{aligned}$$

Now by Markov inequality,

$$\begin{aligned} P[\sqrt{n}\rho_\infty(F_n, F) > M] &\leq M^{-1} E[\sqrt{n}\rho_\infty(F_n, F)] \\ &\leq CM^{-1} \int_0^\infty \exp[-2y^2] dy \rightarrow 0 \quad \text{as } M \rightarrow \infty, \end{aligned}$$

proving that  $\sqrt{n}\rho_\infty(F_n, F) = O_P(1)$  (ie,  $\rho_\infty(F_n, F) = O_P(n^{-1/2})$ ).

The proof of (ii) is longer and we omit it. □

The following theorem provides conditions under which the remainder term

$$R_n = \sqrt{n}[T(F_n) - T(F)] - n^{-1/2} \sum_{i=1}^n T'(X_i; F) = o_P(1).$$

### Theorem 10.6.1.

(i) If  $T$  is  $\rho_\infty$ -Hadamard differentiable at  $F$ , then  $R_n = o_P(1)$ .

(ii) If  $T$  is  $\rho_\infty$ -Fréchet differentiable at  $F$  and if  $\sqrt{n}\rho(F_n, F) = O_P(1)$ , then  $R_n = o_P(1)$ .

*Proof.* Part (i) is proved by advanced techniques, for which we refer to Fernholz [\[57\]](#).

To prove part (ii), note the following:

By definition of  $\rho$ -Fréchet differentiability,  $\rho(F_n, F) \rightarrow 0$  implies

$$\begin{aligned} \frac{T(F_n) - T(F) - L_F(F_n - F)}{\rho(F_n, F)} &= \frac{\sqrt{n}[T(F_n) - T(F)] - \sqrt{n}L_F(F_n - F)}{\sqrt{n}\rho(F_n, F)} \\ &= \frac{R_n}{\sqrt{n}\rho(F_n, F)} \rightarrow 0, \end{aligned}$$

because if  $L_F(F_n - F)$  exists as a  $\rho$ -Fréchet differential, then it is the same as the Gâteaux differential which also exists, so that  $\sqrt{n}L_F(F_n - F) = n^{-1/2} \sum_{i=1}^n T'(X_i; F)$ . Thus for every  $\varepsilon > 0$ , there exists  $\delta > 0$  such that

$$\rho(F_n, F) < \delta \text{ implies } |R_n| < \varepsilon\sqrt{n}\rho(F_n, F) \quad \text{for all large } n.$$

Hence for arbitrary  $\eta > 0$ ,

$$P[|R_n| > \eta] \leq P[\rho(F_n, F) \geq \delta] + P[\sqrt{n}\rho(F_n, F) > \eta/\varepsilon]$$

for large  $n$ , so

$$\limsup_n P[|R_n| > \eta] \leq \limsup_n P[\rho(F_n, F) \geq \delta] + \limsup_n P[\sqrt{n}\rho(F_n, F) > \eta/\varepsilon].$$

Now (i)  $\sqrt{n}\rho(F_n, F) = O_P(1)$  implies  $\rho(F_n, F) = o_P(1)$ ,  $\lim_{n \rightarrow \infty} P[\rho(F_n, F) \geq \delta] = 0$  and (ii)  $\sqrt{n}\rho(F_n, F) = O_P(1)$  also implies  $P[\sqrt{n}\rho(F_n, F) > \eta/\varepsilon]$  can be made arbitrarily small by choosing  $\eta/\varepsilon$  sufficiently large (ie, by making  $\varepsilon$  sufficiently small).

This completes the proof.  $\square$

We now examine the remainder terms of the  $L$ - and  $M$ -estimators.

**Theorem 10.6.2.** *Let  $T(G) = \int xJ(G(x)) dG(x)$  be an  $L$ -functional.*

- (i) *If  $J$  is bounded,  $J(u) = 0$  for all  $u \notin (\alpha, \beta)$  for some  $0 < \alpha < \beta < 1$ , and  $J$  is continuous at  $F(x)$  for almost all  $x$ , then  $T$  is  $\rho_\infty$ -Fréchet differentiable at  $F$  with  $T'(x; F)$  bounded and continuous.*
- (ii) *If  $|J(u) - J(v)| \leq C|u - v|^{p-1}$  for some  $C > 0$  and  $p > 1$ , then  $T$  is  $\rho_{L_p}$ -Fréchet differentiable at  $F$ .*

*Proof.*

- (i) Let  $\{G_j\}$  satisfy  $\lim_{j \rightarrow \infty} \rho(G_j, F) = 0$ . Choose  $c < d$  such that  $F(c) < \alpha < \beta < F(d)$ . Then for  $x \notin [c, d]$ ,  $F(x) \notin (\alpha, \beta)$ , so  $J(F(x)) = 0$  for all  $x \notin [c, d]$ . Since  $\rho_\infty(G_j, F) = 0$ , for sufficiently large  $j$ ,  $G_j(x)$  is also  $\notin (\alpha, \beta)$  for  $x \notin [c, d]$ , so  $J(G_j(x)) = 0$  for all  $x \notin [c, d]$ . Hence for sufficiently large  $j$  and for  $x \notin [c, d]$ ,  $J(u) = 0$  for  $u = F(x)$  and for all  $u$  lying between  $F(x)$  and  $G_j(x)$ , and therefore,  $\int_{F(x)}^{G_j(x)} [J(u) - J(F(x))] du = 0$  and  $W_{G_j}(x) = 0$ .

Hence for large  $j$ ,

$$|R(G_j, F)| = \left| \int_c^d W_{G_j(x)} [G_j(x) - F(x)] dx \right| \leq \rho_\infty(G_j, F) \int_c^d |W_{G_j}(x)| dx.$$

Now use the continuity of  $J$  at  $F(x)$  for almost all  $x$  to see that  $W_{G_j}(x) \rightarrow 0$  a.e., and is bounded since  $J$  is bounded. Hence  $\int_c^d |W_{G_j}(x)| dx \rightarrow 0$  by dominated convergence. Thus  $\lim_{j \rightarrow \infty} |R(G_j, F)| / \rho_\infty(G_j, F) = 0$  showing that  $T$  is  $\rho_\infty$ -Fréchet differentiable at

$F$ . Again, for  $y \notin [c, d]$ ,  $F(y) \notin [\alpha, \beta]$  so  $J(F(y)) = 0$ . Thus

$$T'(x; F) = - \int_c^d [\delta_x(y) - F(y)] J(F(y)) dy$$

is bounded and continuous. This proves (i).

(ii) Let  $\{G_j\}$  satisfy  $\lim_{j \rightarrow \infty} \rho_{L_p}(G_j, F) = 0$  and note that

$$\begin{aligned} |R(G_j, F)| &= \left| \int \left[ \int_{F(x)}^{G_j(x)} \{J(u) - J(F(x))\} du \right] dx \right| \\ &\leq \int \left[ \int_{F(x)}^{G_j(x)} |J(u) - J(F(x))| du \right] dx \\ &\leq \int \left[ \int_{F(x)}^{G_j(x)} C|u - F(x)^{p-1}| du \right] dx \\ &\leq C \int |G_j(x) - F(x)|^p dx = C \left\{ \rho_{L_p}(G_j, F) \right\}^p. \end{aligned}$$

Hence

$$\lim_{j \rightarrow \infty} \frac{|R(G_j, F)|}{\rho_{L_p}(G_j, F)} \leq C \lim_{j \rightarrow \infty} \left\{ \rho_{L_p}(G_j, F) \right\}^{p-1} = 0 \quad \text{for } p > 1,$$

showing that  $T$  is  $\rho_{L_p}$ -Fréchet differentiable at  $F$ .  $\square$

**Corollary 10.6.1.** *In the expansion*

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n$$

of the  $L$ -estimator  $T(F_n)$  of  $T(f) = \int xJ(F(x)) dx$ , the remainder term  $R_n = o_P(1)$  if the function  $J$  satisfies the conditions of Theorem 10.6.2(i) or (ii).

*Proof.* Under the condition of Theorem 10.6.2(i),  $T$  is  $\rho_\infty$ -Fréchet differentiable at  $F$  and by Lemma 10.6.1(i),  $\rho_\infty(F_n, F) = O_P(n^{-1/2})$ . Therefore, Theorem 10.6.1(ii) applies, proving  $R_n = o_P(1)$ .

Under the condition of Theorem 10.6.2(ii) with  $p > 1$ ,  $T$  is  $\rho_{L_p}$ -Fréchet differentiable and by Lemma 10.6.1(ii),  $\rho_{L_p}(F_n, F) = O_P(n^{-1/2})$  provided that  $\int [F(x)(1 - F(x))^{p/2}] dx < \infty$  also holds for the case  $1 < p < 2$ . Hence Theorem 10.6.1(ii) applies, proving  $R_n = o_P(1)$ .  $\square$

**Theorem 10.6.3.** *Let  $T$  be an  $M$ -functional with score function  $\Psi(x, t)$  which is bounded and continuous, and suppose that the function  $\lambda_F(t) = \int \Psi(x, t) dF(x)$  is continuously differentiable at  $T(F)$  while  $\lambda'_F(T(F)) \neq 0$ . Then  $T$  is  $\rho_\infty$ -Hadamard differentiable at  $F$ .*

*Proof.* Let  $G_j = F + t_j \Delta_j$  where  $\|\Delta_j - \Delta\|_\infty \rightarrow 0$  on  $\mathcal{D} = \{c(F_1 - F_2), c > 0\}$  and  $t_j \rightarrow 0$ . We argue through the following steps:

- (i) Since the  $\rho_\infty$ -Hadamard differential, if it exists, is the same as the Gâteaux differential,

$$L_F(G_j - F) = \int T'(x; F) d[G_j(x) - F(x)] = - \int \frac{\Psi(x, T(F))}{\lambda'_F(T(F))} d[G_j(x) - F(x)],$$

using the formula for  $T'(x; F)$  given in Eq. (11). It will, therefore, be enough to show that

$$\begin{aligned} & \lim_{j \rightarrow \infty} \frac{T(G_j) - T(F) + \{\lambda'_F(T(F))\}^{-1} \int \Psi(x, T(F)) d[G_j(x) - F(x)]}{\|t_j \Delta_j\|_\infty} \\ &= \lim_{j \rightarrow \infty} R_j / \|t_j \Delta_j\|_\infty, \end{aligned} \quad (15)$$

where  $R_j$  is the numerator of Eq. (15).

- (ii) Since  $\lambda_F(T(F)) = \lambda_{G_j}(T(G_j)) = 0$  by (12),

$$\begin{aligned} \lambda_F(T(G_j)) - \lambda_F(T(F)) &= \lambda_F(T(G_j)) - \lambda_{G_j}(T(G_j)) \\ &= - \int \Psi(x, T(G_j)) d[G_j(x) - F(x)] \\ &= - \int \Psi(x, T(G_j)) d[t_j \Delta_j(x)] \rightarrow 0, \end{aligned}$$

because  $t_j \rightarrow 0$ ,  $\|\Delta_j - \Delta\|_\infty \rightarrow 0$  and  $\Psi$  is bounded.

- (iii) Since  $\lambda'_F(T(F)) \neq 0$ , the inverse function  $\lambda_F^{-1}(\cdot)$  of  $\lambda_F(\cdot)$  exists and is continuous in a neighborhood of  $\lambda_F(T(F)) = 0$ . Hence  $\lambda_F(T(G_j)) - \lambda_F(T(F)) \rightarrow 0$  implies

$$T(G_j) - T(F) = \lambda_F^{-1}(\lambda_F(T(G_j))) - \lambda_F^{-1}(\lambda_F(T(F))) \rightarrow 0.$$

- (iv) Use (ii) to write

$$T(G_j) - T(F) = - \frac{T(G_j) - T(F)}{\lambda_F(T(G_j)) - \lambda_F(T(F))} \int \Psi(x, T(G_j)) d[t_j \Delta_j(x)].$$

- (v) From (i) and (iv),

$$\begin{aligned} R_j &= \frac{1}{\lambda'_F(T(F))} \int \Psi(x, T(F)) d[t_j \Delta_j(x)] \\ &\quad - \frac{T(G_j) - T(F)}{\lambda_F(T(G_j)) - \lambda_F(T(F))} \int \Psi(x, T(G_j)) d[t_j \Delta_j(x)] \\ &= t_j \left[ \left\{ \lambda'_F(T(F)) \right\}^{-1} - \left\{ \frac{\lambda_F(T(G_j)) - \lambda_F(T(F))}{T(G_j) - T(F)} \right\}^{-1} \int \Psi(x, T(F)) d\Delta_j(x) \right] \end{aligned}$$

$$\begin{aligned}
& - t_j \frac{T(G_j) - T(F)}{\lambda_F(T(G_j)) - \lambda_F(T(F))} \int [\Psi(x, T(G_j)) - \Psi(x, T(F))] d\Delta_j(x) \\
& := t_j(R_{1j} + R_{2j}).
\end{aligned}$$

(vi) By (iii),  $T(G_j) - T(F) \rightarrow 0$ , so

$$\left\{ \frac{\lambda_F(T(G_j)) - \lambda_F(T(F))}{T(G_j) - T(F)} \right\}^{-1} - \{\lambda'_F(T(F))\}^{-1} \rightarrow 0,$$

which together with boundedness of  $\Psi$  and  $\|\Delta_j - \Delta\|_\infty \rightarrow 0$  implies  $R_{1j} \rightarrow 0$ .

Also  $T(G_j) - T(F) \rightarrow 0$  implies

$$\begin{aligned}
& \frac{T(G_j) - T(F)}{\lambda_F(T(G_j)) - \lambda_F(T(F))} \rightarrow \{\lambda'_F(T(F))\}^{-1} \text{ and} \\
& \int [\Psi(x, T(G_j)) - \Psi(x, T(F))] d\Delta_j(x) \rightarrow 0,
\end{aligned}$$

because  $\Psi$  is bounded and continuous, and  $\|\Delta_j - \Delta\|_\infty \rightarrow 0$ . Thus  $R_{2j} \rightarrow 0$ . Putting all this together, we have

$$\lim_{j \rightarrow \infty} \frac{R_j}{\|t_j \Delta_j\|_\infty} = \lim_{j \rightarrow \infty} \frac{t_j(R_{1j} + R_{2j})}{\|t_j \Delta_j\|_\infty} = \lim_{j \rightarrow \infty} \frac{R_{1j} + R_{2j}}{\|\Delta_j\|_\infty} = 0,$$

showing that Eq. (15) holds.  $\square$

**Corollary 10.6.2.** *If  $T(F)$  is an M-functional with score function  $\Psi(x, t)$  satisfying the conditions of Theorem 10.6.3 and if  $T(F_n)$  is the corresponding M-estimator, then in the expansion*

$$\sqrt{n}[T(F_n) - T(F)] = n^{-1/2} \sum_{i=1}^n T'(X_i; F) + R_n,$$

the remainder term  $R_n$  is  $o_P(1)$ .

*Proof.* Under the conditions of Theorem 10.6.3,  $T$  is  $\rho_\infty$ -Hadamard differentiable at  $F$ , so Theorem 10.6.1(i) applies.  $\square$

## 10.7 The Jackknife and the Bootstrap

In this section we will briefly deal with two well-known resampling methods — the jackknife and the bootstrap. These methods can be used to obtain approximate bias, variance, and distribution of estimates without having to obtain their analytic expressions which may be quite complicated in many cases (see Efron and Tibshirani [59]). At the beginning of this chapter, some examples of statistical functional were given, and in

[Examples 10.2.1–10.2.3](#), the first derivatives (influence functions) were explicitly obtained for three functionals. It is important to add that there are many more useful functionals for the univariate case as in [Exercises 10.5–10.11](#) or in the multivariate case such as the correlation coefficient, multiple correlation, etc. In this section, the discussion will be informal, starting with the concept of second derivative of a functional. Extending the definition given in [Section 10.2](#) (for the first-order expansion of a statistical functional),  $T''(x_1, x_2; F)$  is the second derivative if

$$\begin{aligned} T(F + t\Delta) - T(F) - t \int T'(x; F) d\Delta(x) - t^2 \int T''(x_1, x_2; F) d\Delta(x_1) d\Delta(x_2) \\ = o(t^2), \end{aligned}$$

as  $t \downarrow 0$ . We ignore the issue of whether the convergence is uniform in  $\Delta$ . This has been discussed in the previous sections. We now examine the bias and the variance of the estimate  $T(F_n)$  of  $T(F)$ , where  $F_n$  is the empirical cdf on the basis of iid sample  $X_1, \dots, X_n$ . If we write  $\Delta_n = \sqrt{n}(F_n - F)$  and  $t = n^{-1/2}$ , under appropriate conditions we have

$$\begin{aligned} T(F_n) &= T(F) + n^{-1/2} \int T'(x; F) d\Delta_n(x) \\ &\quad + n^{-1} \int T''(x_1, x_2; F) d\Delta_n(x_1) d\Delta_n(x_2) + R_n \\ &:= T(F) + L_n(F) + Q_n(F) + R_n, \end{aligned} \tag{16}$$

where we assume that the remainder term is  $o_P(1/n)$ . However, we should point out that, usually for many functionals,  $R_n = O_P(n^{-3/2})$ .

We now discuss the concepts of asymptotic bias (ABias) and asymptotic variance (AVar) of a statistical functional. In all our subsequent discussions in this section, *we ignore the remainder term  $R_n$  and assume that  $T(F_n) = T(F) + L_n(F) + Q_n(F)$ , and that  $E[\{T'(X; F)\}^2]$ ,  $E[\{T''(X, X; F)\}^2]$ , and  $E[\{T''(X_1, X_2; F)\}^2]$  are finite.*

Note that  $L_n(F)$  has mean zero and variance

$$\text{Var}[L_n(F)] = n^{-1} \text{Var}[T'(X; F)].$$

A tedious calculation, whose justification will be given later, shows

$$\text{Var}[L_n(F) + Q_n(F)] = \text{Var}[L_n(F)] + O(n^{-2}).$$

The expected value of  $Q_n(F)$  is

$$\begin{aligned} E[Q_n(F)] &= n^{-1} [E\{T''(X, X; F)\} - E\{T''(X_1, X_2; F)\}], \text{ and hence} \\ E[L_n(F) + Q_n(F)] &= E[Q_n(F)] = n^{-1} [E\{T''(X, X; F)\} - E\{T''(X_1, X_2; F)\}]. \end{aligned}$$

We define the asymptotic bias and asymptotic variance of  $T(F_n)$  as

$$\begin{aligned} \text{ABias}(T(F_n)) &= E[Q_n(F)] \\ &= n^{-1} [E\{T''(X, X; F)\} - E\{T''(X_1, X_2; F)\}] := n^{-1} b(F), \end{aligned} \tag{17}$$

$$\text{AVar}[T(F_n)] = \text{Var}[L_n(F)] = n^{-1} \text{Var}[T'(X; F)] := n^{-1} v(F). \quad (18)$$

In some cases, it is possible to find explicit expressions for ABias and AVar, and in such cases, we may simply replace  $b(F)$  and  $v(F)$  by  $b(F_n)$  and  $v(F_n)$ . However, in many others, explicit expressions for ABias and AVar are quite complicated, and it is useful to have simple sample-based methods for estimating these quantities. We now discuss estimation of these two quantities: asymptotic bias and asymptotic variance.

*Remark 10.7.1.* It is important to note that  $T(F_n)$  may not have finite mean and variance in many cases. For instance, let  $X_1, \dots, X_n$  be iid discrete random variables taking values in  $\mathbb{N} = \{0, 1, 2, \dots\}$ . Let  $\mu(F) = E[X]$  and  $\sigma^2(F) = \text{Var}[X]$ , and we are interested in estimating  $T(F) = \log(\mu(F))$  or the coefficient-of-variation  $T(F) = \sigma(F)/\mu(F)$ . Note that if  $P[X = 0] > 0$ , then  $T(F_n)$  in each of these cases does not have finite mean and the variance is not defined. However, in statistical applications, the issues of interest are estimation of  $T(F)$  and construction of its confidence interval. Hence if  $P[T(F_n) = \infty] \rightarrow 0$  as  $n \rightarrow \infty$ , we may bypass the problem of estimating the actual bias and variance (may not exist) by assuming that the expression for  $T(F_n)$  is given by the expression on the right-hand side of Eq. (16) without the remainder term  $R_n$ .

### 10.7.1 Estimation of Asymptotic Bias and Asymptotic Variance

Let  $F_{ni}$  be the empirical cdf of  $F$  on the basis of  $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n$  (ie,  $F_{ni}$  is based on  $n - 1$  observations, deleting  $X_i$ ). Let  $T(F_{ni})$  be the estimate of  $T(F)$  based on  $F_{ni}$ . Tukey's pseudo values are defined to be

$$V_i = nT(F_n) - (n - 1)T(F_{ni}), \quad i = 1, \dots, n.$$

These pseudo values are *approximately iid* with mean  $T(F)$  and variance  $v(F)/n$ . As a matter of fact, the following turn out to be true

$$\begin{aligned} E[\bar{V}] &= T(F), \\ E\left[(n-1)^{-1} \sum_{i=1}^n (V_i - \bar{V})^2\right] &= v(F) + O(n^{-1}), \end{aligned}$$

where  $\bar{V} = n^{-1} \sum_{i=1}^n V_i$ . Since  $E[T(F_n) - \bar{V}] = E[T(F_n)] - T(F)$ , the Jackknife estimates of ABias( $T(F_n)$ ) is given by

$$\text{ABias}^{(J)}(T(F_n)) = T(F_n) - \bar{V} = (n-1)[\bar{T}_n - T(F_n)],$$

where  $\bar{T}_n = n^{-1} \sum_{i=1}^n T(F_{ni})$ .

The Jackknife estimate of AVar( $T(F_n)$ ) is

$$\text{AVar}^{(J)}(T(F_n)) = [n(n-1)]^{-1} \sum_{i=1}^n (V_i - \bar{V})^2 = \frac{n-1}{n} \sum_{i=1}^n (T(F_{ni}) - \bar{T}_n)^2.$$

The bootstrap method is conceptually simple: it basically seeks to replace  $b_n(F)$  by its empirical estimate  $b_n(F_n)$ . However, it does so without having to obtain any analytic ex-

pression for  $b_n(F)$ . The method described here is known as the “nonparametric” bootstrap. Let  $X_1^*, \dots, X_n^*$  be iid with cdf  $F_n$ , then the bootstrap estimate of ABias is given by

$$\text{ABias}^{(B)}(T(F_n)) = E[\{T(F_n^*) - T(F_n)\}|X_1, \dots, X_n],$$

that is, the (conditional) expectation is taken over the random sample  $X_1^*, \dots, X_n^*$ . In practice, it is calculated as follows. Draw a random sample of size  $n$  with replacement from the data  $\{X_1, \dots, X_n\}$  and repeat this  $N$  times. Let  $F_{n(t)}^*$  be the empirical cdf on the basis of the  $t$ th sample  $\{X_{1(t)}^*, \dots, X_{n(t)}^*\}$ , and let  $T(F_{n(t)}^*)$  be the estimate of  $T(F)$  based on  $F_{n(t)}^*$ . Then, one calculates the quantity

$$N^{-1} \sum_{t=1}^N [T(F_{n(t)}^*) - T(F_n)],$$

which, by the weak law of large numbers, converges to  $\text{ABias}^{(B)}(T(F_n))$  as  $N \rightarrow \infty$ .

The bootstrap estimate of the asymptotic variance is

$$\text{AVar}^{(B)}(T(F_n)) = \text{Var}[T(F_n^*)|X_1, \dots, X_n],$$

where the conditional variance is over the bootstrap sample  $(X_1^*, \dots, X_n^*)$ .

The bootstrap procedure may also be used to obtain an estimate of the sampling distribution  $Q_n(z) = P[\sqrt{n}\{T(F_n) - T(F)\} \leq z]$ ,  $z$  real. Even though such sampling distributions are approximately normal under appropriate conditions when the sample size  $n$  is large, one may nonetheless use the bootstrap method in such cases. A bootstrap estimate of  $Q_n(z)$  is given by  $Q_n^{(B)}(z) = P[\sqrt{n}\{T(F_n^*) - T(F_n)\} \leq z|X_1, \dots, X_n]$ , where the (conditional) probability is over the bootstrap sample  $(X_1^*, \dots, X_n^*)$ . In general if there is a functional of the form  $T(F_n, F)$  and one wishes to estimate  $E[T(F_n, F)]$ , then its bootstrap estimate is  $E[T(F_n^*, F_n)|X_1, \dots, X_n]$ .

*Remark 10.7.2.* It is important to point out that the jackknife and bootstrap procedures may not always work. For instance, the jackknife estimate of bias cannot provide consistent estimates for the ABias and AVar when estimating a quantile  $F^{-1}(p)$ . Success of the jackknife method depends on the smoothness of the functional (ie, on the validity of the expansion given in Eq. (16)). Bootstrap works well for quantile estimation as long as  $p$  is away from 0 or 1. However, it cannot provide consistent estimates for ABias and AVar when estimating extreme quantiles (ie, when  $p$  is close to 0 or 1).

### 10.7.2 Heuristic Justification for the Jackknife and the Bootstrap

Let us first justify the validity of the bootstrap estimate. In the arguments given here we denote  $n^{-1/2} \Delta_n = F_n - F$  by  $D_n$  and  $F_n^* - F_n$  by  $D_n^*$ . We can simplify notations by writing  $\int T'(x; F) dD_n(x)$  by  $L_F(D_n)$  and  $\int T''(x_1, x_2; F) dD_n(x_1) dD_n(x_2)$  by  $Q_F(D_n)$ . If we expand  $T(F_n^*)$  about  $F_n$ , then we have

$$T(F_n^*) - T(F_n) = L_{F_n}(D_n^*) + Q_{F_n}(D_n^*).$$

Conditional expectation (given  $X_1, \dots, X_n$ ) of the first term on the right-hand side of the last expression equals zero, and the conditional expectation of the second term is  $n^{-1}b_n(F_n)$ . The result follows once we note that  $|b_n(F_n) - b_n(F)| \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . The justification for the bootstrap estimate of the variance is similar.

Let us now look at the jackknife estimate. The validity of the bias estimate is fairly easy to establish since

$$\mathbb{E}[V_i] = \mathbb{E}[nT(F_n) - (n-1)T(F_{ni})] = b(F) - b(F) = 0.$$

Thus  $\mathbb{E}[\bar{V}] = 0$ .

Justification for the validity of the jackknife estimate of the variance requires a bit of work. Simple algebra tells us

$$\begin{aligned} \sum_{i=1}^n (V_i - \bar{V})^2 &= (2n)^{-1} \sum_{1 \leq i \neq i' \leq n} (V_i - V_{i'})^2 \\ &= (2n)^{-1}(n-1)^2 \sum_{1 \leq i \neq i' \leq n} [T(F_{ni}) - T(F_{ni'})]^2. \end{aligned} \quad (19)$$

We will work with  $T(F_{ni}) - T(F_{ni'})$  for the proof.

For notational simplicity, we suppress  $F$  in the notations of  $T'$  and  $T''$ . We also assume that  $T''$  is symmetric in its arguments. Since if it is not, we can symmetrize it by using  $[T''(x_1, x_2) + T''(x_2, x_1)]/2$ , and such symmetrization does not alter the value of  $\int T''(x_1, x_2) dH(x_1) dH(x_2)$ , where  $H$  is a cdf or a difference of two cdfs. For any  $x, x_1$ , and  $x_2$ , define

$$\begin{aligned} T_1(x) &= T'(x) - \mathbb{E}[T'(X)], \\ T_2(x_1, x_2) &= T''(x_1, x_2) - \mathbb{E}[T''(x_1, X_2)] - \mathbb{E}[T''(X_1, x_2)] + \mathbb{E}[T''(X_1, X_2)]. \end{aligned}$$

Then

$$\begin{aligned} L_n(F) &= \int T'(x) dD_n(x) = \int T_1(x) dF_n(x) = n^{-1} \sum_{j=1}^n T_1(X_j), \\ Q_n(F) &= \int T''(x_1, x_2) dD_n(x_1) dD_n(x_2) \\ &= \int T_2(x_1, x_2) dF_n(x_1) dF_n(x_2) = n^{-2} \sum_{1 \leq j, k \leq n} T_2(X_j, X_k). \end{aligned}$$

Thus

$$T(F_n) = T(F) + n^{-1} \sum_{j=1}^n T_1(X_j) + n^{-2} \sum_{1 \leq j, k \leq n} T_2(X_j, X_k).$$

The two sums in the last expression are not uncorrelated. A simple argument can be used to create two uncorrelated sums. Note that  $b(F) = \mathbb{E}[T_2(X, X)]$  and let

$$\begin{aligned} U_i &= T_1(X_i) + n^{-1}[T_2(X_i, X_i) - b(F)], \\ W_{j,k} &= T_2(X_j, X_k). \end{aligned}$$

Then  $T(F_n)$  can be rewritten as

$$T(F_n) = T(F) + b(F)/n + n^{-1} \sum_{j=1}^n U_i + n^{-2} \sum_{1 \leq j \neq k \leq n} W_{j,k}. \quad (20)$$

It is easy to see that  $U_i$  and  $W_{j,k}$ ,  $j \neq k$ , have zero means,  $\text{Cov}[U_i, W_{j,k}] = 0$  whenever  $j \neq k$ , and  $\text{Cov}[W_{j_1,k_1}, W_{j_2,k_2}] = 0$  except when  $(j_1, k_1) = (j_2, k_2)$  or  $(j_1, k_1) = (k_2, j_2)$ . We therefore have

$$\begin{aligned} \mathbb{E}[T(F_n)] &= T(F) + b(F)/n, \\ \text{Var}[T(F_n)] &= \text{Var}\left[n^{-1} \sum_{j=1}^n U_i\right] + \text{Var}\left[n^{-2} \sum_{1 \leq j \neq k \leq n} W_{j,k}\right] \\ &= n^{-1} \text{Var}[U_1] + O(n^{-2}) \\ &= n^{-1} \{\text{Var}[T_1(X)] + O(n^{-1})\} + O(n^{-2}) \\ &= n^{-1} v(F) + O(n^{-2}). \end{aligned}$$

Applying the equality in Eq. (20) for  $T(F_{ni})$  and  $T(F_{ni'})$  for  $i \neq i'$ , we have

$$T(F_{ni}) - T(F_{ni'}) = (n-1)^{-1} [U_{i'} - U_i] + 2(n-1)^{-2} \sum_{j \neq i, j \neq i'} [W_{i',j} - W_{i,j}].$$

Thus

$$\begin{aligned} \mathbb{E} \left\{ \sum_{1 \leq i \neq i' \leq n} [T(F_{ni}) - T(F_{ni'})]^2 \right\} &= \sum_{i \neq i'} \mathbb{E} \left\{ (n-1)^{-1} [U_{i'} - U_i] \right\}^2 \\ &\quad + \sum_{i \neq i'} \mathbb{E} \left\{ 2(n-1)^{-2} \sum_{j \neq i, j \neq i'} [W_{i',j} - W_{i,j}] \right\}^2 \\ &= 2n(n-1)^{-1} \text{Var}[U_1] + O(n^{-1}) \\ &= 2n(n-1)^{-1} \{\text{Var}[T_1(X)] + O(n^{-1})\} + O(n^{-1}) \\ &= 2n(n-1)^{-1} v(F) + O(n^{-1}). \end{aligned}$$

From Eq. (19), we thus have

$$\mathbb{E} \left[ \sum_{i=1}^n (V_i - \bar{V})^2 \right] = (n-1) \text{Var}[T_1(X)] + O(1) = (n-1)v(F) + O(1).$$

This shows that  $\mathbb{E}[A\text{Var}^{(J)}(T(F_n))] = n^{-1}v(F) + O(n^{-2})$ .

## Exercises

In the following problems,  $X_1, \dots, X_n$  is a random sample from a population with cdf  $F$  and pdf  $f$ ,  $\bar{X}_n$  and  $s_n^2$  are the sample mean and sample variance,  $X_{n:1} < \dots < X_{n:n}$  are the

order statistics, and  $Q_{1n} = X_{n:[n/4]}$ ,  $M_n = X_{n:[n/2]}$ , and  $Q_{3n} = X_{n:[3n/4]}$  are, respectively, the sample first quartile, the sample median, and the sample third quartile.

- 10.1. Let  $F$  be the cdf of  $N(\mu, \sigma^2)$ . Consider two tests with critical regions  $T_{1n} \geq c_{1n}$  and  $T_{2n} \geq c_{2n}$  for testing  $H_0: \mu = 0$  vs  $H_1: \mu > 0$ , where  $T_{1n} = \sqrt{n}\bar{X}_n/s_n$ ,  $T_{2n} = Q_{1n} + Q_{3n}$ , and  $c_{1n}, c_{2n}$  are chosen so as to control the Type I error probability at a given level  $\alpha$ , approximately for large  $n$ . For alternatives of the order  $1/\sqrt{n}$ , calculate the Pitman ARE of the test based on  $T_{2n}$  with respect to the one based on  $T_{1n}$ .
- 10.2. If  $F$  is the cdf of  $N(\mu, \sigma^2)$ , find the asymptotic joint distribution of  $(\bar{X}_n, M_n)$ , properly normalized.
- 10.3. Let  $\xi_1$  and  $\xi_3$  denote the first and the third quartiles of  $F$  (ie,  $\xi_1 = F^{-1}(1/4)$  and  $\xi_3 = F^{-1}(3/4)$ ).
  - (a) Write down the Bahadur representation of  $Q_{1n}$  and  $Q_{3n}$ , stating conditions on  $F$  for their validity.
  - (b) Let  $F$  be the cdf of  $N(\mu, \sigma^2)$ . Find the asymptotic distribution of the inter-quartile range  $D_n = Q_{3n} - Q_{1n}$ .
  - (c) Find a constant  $c$  such that  $T_n = cD_n$  is a consistent estimator of  $\sigma$ . What can you say about the asymptotic efficiency of  $T_n$ ? [If  $\phi$  is the pdf and  $\Phi$  is cdf of  $N(0, 1)$ , then  $\Phi^{-1}(3/4) = 0.6745$  and  $\phi(\Phi^{-1}(3/4)) = 0.3178$ .]
- 10.4. Suppose the pdf involves a parameter  $\theta$  and

$$\begin{aligned} f(x, \theta) &= (1/2 - \theta)e^x I_{(-\infty, 0)}(x) + (1/2)e^{-x/(1+2\theta)} \\ &= \begin{cases} (1/2 - \theta)e^{-|x|} & x < 0 \\ (1/2 + \theta)\frac{1}{1+2\theta}e^{-|x|/(1+2\theta)} & x \geq 0. \end{cases} \end{aligned}$$

For  $\theta = 0$ , the pdf  $f(x, 0) = (1/2)e^{-|x|}$ ,  $-\infty < x < \infty$  is symmetric with mean = median = 0. For  $\theta \neq 0$ , none of this is true. To test  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  we can, therefore, test for symmetry on the basis of the mean or the median. Another possibility is to test whether the MLE  $\hat{\theta}_n$  of  $\theta$  exhibits distributional property that it should have under  $H_0$ . With this in mind, consider the following test statistics:

$$\begin{aligned} T_{0n} &= \hat{\theta}_n, \quad T_{1n} = \bar{X}_n, \quad T_{2n} = M_n, \\ T_{3n} &= Q_{3n} - 2M_n + Q_{1n} \text{ (a measure of asymmetry).} \end{aligned}$$

- (a) Find the asymptotic distributions of  $T_{0n}$ ,  $T_{1n}$ ,  $T_{2n}$ , and  $T_{3n}$  (properly normalized).
- (b) Consider tests with critical regions  $T_{jn} \geq c_{jn}$ ,  $j = 0, 1, 2, 3$ , each  $c_{jn}$  being chosen to control the Type I error probabilities at a given level  $\alpha$  approximately, for large  $n$ . For alternative  $\theta > 0$  of the order  $1/\sqrt{n}$ , calculate the Pitman AREs of the tests based on  $T_{1n}$ ,  $T_{2n}$ ,  $T_{3n}$  with respect to the one based on  $T_{0n}$ .

- 10.5.** (a) Let  $T_1$  and  $T_2$  be statistical functionals of a family of cdf's  $\mathcal{F}_0$  having finite third moments. Express the influence function  $T'(x; F)$  of  $T = T_1/T_2$  in terms of  $T_1(F)$ ,  $T_2(F)$  and their influence functions  $T'_1(x; F)$  and  $T'_2(x; F)$ .
- (b) Let  $T(F) = E_F[(X - \xi(F))^3]/\{E_F[(X - \xi(F))^2]\}^{3/2}$ , where  $\xi(F) = E_F[X]$ . Use the result of (a) and the formula for the influence function of  $E_F[(X - \xi(F))^k]$  given in Section 10.2 to derive the influence function of this  $T(F)$ .
- (c) Verify for this  $T(F)$  that  $E_F[T'(X; F)] = 0$ ,  $R_n = \sqrt{n}[T(F_n) - T(F)] - n^{-1/2} \sum_{i=1}^n T'(X_i; F) = o_P(1)$ , and then find the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$  for arbitrary  $F \in \mathcal{F}_0$  and when  $F$  is the cdf of  $N(\mu, \sigma^2)$ .
- 10.6.** Let  $\mathcal{F}_0$  be the set of all cdf's on  $\mathbb{R}$  with finite mean.
- (a) Find  $T(F)$  on  $\mathcal{F}_0$  such that  $T(F_n) = n^{-1} \sum_{i=1}^n |X_i - \bar{X}_n|$ .
- (b) Find the influence function of  $T(F)$ .
- (c) Specializing to the case of symmetric  $F$ , indicate how you would show that  $R_n = \sqrt{n}[T(F_n) - T(F)] - n^{-1/2} \sum_{i=1}^n T'(X_i; F) = o_P(1)$ .
- 10.7.** Let  $T(F) = E_F[|X - X'|]$ , where  $X, X'$  are iid with cdf  $F$ .
- (a) Find the influence function of  $T(F)$ .
- (b) Write down the expression for  $T(F_n)$  and derive the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$ .
- 10.8.** Let  $T(F) = (\int x dF(x))^{1/2} = \xi(F)^{1/2}$  for  $F$  having finite mean  $\xi(F) > 0$ .
- (a) Find the influence function of  $T'(x; F)$ .
- (b) Verify that  $E_F[T'(X; F)] = 0$ ,  $R_n = \sqrt{n}[T(F_n) - T(F)] - n^{-1/2} \sum_{i=1}^n T'(X_i; F) = o_P(1)$ , and then find the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$ .
- 10.9.** Let  $T(F) = [\int (x - \xi(F))^2 dF(x)]^{1/2}$  for a cdf with mean  $E_F[X] = \xi(F)$  and finite second moment. Find the influence function of  $T(F)$  and the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$ .
- 10.10.** For cdf's on  $\mathbb{R}$  with nonzero mean  $\xi(F)$ , define  $T(F) = 1/\xi(F)$ . Find the influence function of  $T(F)$  and the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$ .
- 10.11.** For cdf's on  $\mathbb{R}$  with positive mean  $\xi(F)$  and finite second moment, let  $T(F) = [\int (x - \xi(F))^2 dF(x)]^{1/2}/\xi(F)$  denote the coefficient of variation. Find the influence function of  $T(F)$  using the results of Exercises 10.5(a), 10.9, and 10.10, and then find the asymptotic distribution of  $\sqrt{n}[T(F_n) - T(F)]$ .
- 10.12.** The pdf  $f(x, \theta)$  of  $X$  is of the form

$$f(x, \theta) = \begin{cases} 1/4 & |x - \theta| \leq 1 \\ (1/4) \exp[-|x - \theta| + 1] & |x - \theta| > 1, \end{cases}$$

which is uniform in the middle and double-exponential in the tails.

- (a) Find formulas for the  $\alpha$ -quantile and the  $(1 - \alpha)$ -quantile of this distribution for  $0 < \alpha < 1/4$  and for  $1/4 < \alpha < 1/2$ .

- (b) Let  $\bar{X}_{n(\alpha)} = (n - 2[\alpha n])^{-1} \sum_{i=[\alpha n]+1}^{n-[\alpha n]} X_{n:i}$  denote the  $\alpha$ -trimmed mean based on  $X_1, \dots, X_n$ . Find the asymptotic distribution of  $\sqrt{n}[\bar{X}_{n(\alpha)} - \theta]$ .
- (c) Plot the asymptotic variance  $\sigma_1^2(\alpha)$  obtained in (b) for  $\alpha = 0.1, 0.2, 0.3, 0.4$  to determine the best choice of trimming among these.
- (d) Let  $M_n(k)$  denote the  $M$ -estimator of  $\theta$  obtained by solving  $\sum_{i=1}^n \Psi(X_i - t) = 0$ , using the Huber function  $\Psi(x) = xI(|x| \leq k) + k \operatorname{sign}(x)I(|x| > k)$ . Find the asymptotic distribution of  $\sqrt{n}[M_n(k) - \theta]$ .
- (e) Plot the asymptotic variance  $\sigma_2^2(k)$  obtained in (d) for suitably chosen values of  $k$  (starting with  $k = 0.5, 1.0, 1.5, 2.0$ ) to determine as good a choice of  $k$  you can.
- (f) Compare the performances of  $\bar{X}_{n(\alpha)}$  with  $\alpha$  chosen in (c) and  $M_n(k)$  with  $k$  chosen in (e).

In the following two problems, the population cdf and pdf are, respectively,  $F_\theta(x) = F(x - \theta)$  and  $f_\theta(x) = f(x - \theta)$ , where  $f$  is symmetric about 0.

- 10.13.** (a) Suppose that  $T$  is an  $L$ -functional with score function  $J$ . Show that if  $J$  is symmetric about 1/2 with  $\int_0^1 J(u) du = 1$ , then  $T(F_\theta) = \theta$ .
- (b) Let  $\bar{X}_{n(\alpha)}$  denote the  $\alpha$ -trimmed mean. Find the asymptotic variance of  $\sqrt{n}[\bar{X}_{n(\alpha)} - \theta]$  when  $\alpha = 1/4$  and  $f(x)$  is the pdf of the standard Cauchy distribution.
- 10.14.** Let  $T(F)$  be the solution of the equation  $\int \Psi(x - t) dF(x) = 0$  and let  $\hat{\theta}_n$  be the corresponding  $M$ -estimator. Make the following assumptions as needed:
- (i) differentiation under the integral sign is valid,
  - (ii)  $xf(x)$  is of bounded variation,
  - (iii)  $\int |\Psi(x)|f'(x) dx < \infty$  and  $\neq 0$ ,
  - (iv)  $\int \Psi^2(x)f(x) dx < \infty$ .
- (a) Show that if  $\Psi(-x) = -\Psi(x)$  for all  $x$ , then  $T(F_\theta) = \theta$ .
- (b) Let  $0 < I_f = \int \{f'(x)/f(x)\}^2 f(x) dx < \infty$  denote the Fisher-information of the location family  $\{f(x - \theta), \theta \in \mathbb{R}\}$  and let  $\sigma^2(F, \Psi)$  denote the asymptotic variance of  $\sqrt{n}(\hat{\theta}_n - \theta)$ . Show that  $\sigma^2(F, \Psi) \geq I_f^{-1}$ . [Hint: Use integration by parts to rewrite the denominator of  $\sigma^2(F, \Psi)$  and then apply Cauchy-Schwarz inequality.]
- (c) Calculate  $\sigma^2(F, \Psi)$  when  $f$  is the pdf of the standard Cauchy distribution and  $\Psi(x) = -xI_{[-c, c]}(x)$ . [Hint: Show that  $\int_0^c (1+x^2)^{-1} dx = \arctan c$  and  $\int_0^c x^2(1+x^2)^{-1} dx = c - \arctan c$ .]

- 10.15.** The sample mean in a random sample from a Cauchy distribution can be stabilized by excluding just one- or two-order statistics at each end. Let  $\bar{X}_{n,r}^* = (n - 2r)^{-1} \sum_{i=r+1}^{n-r} X_{n:i}$ . Show that  $\bar{X}_{n,r}^*$  is an unbiased estimator of the Cauchy median  $\theta$  if  $r \geq 1$  and that  $\bar{X}_{n,r}^*$  has finite variance if  $r \geq 2$ .

- 10.16.** Calculate the ARE of the  $M$ -estimator with score function  $\Psi(x) = -xI_{[-c,c]}(x)$  of the median  $\theta$  of a Cauchy distribution, comparing  $\sigma^2(F, \Psi)$  with  $I_f^{-1}$ . Use the results of Exercise 10.14(c) with  $c = 1$ .
- 10.17.** Let  $F_\theta(x) = F(x - \theta)$  be a cdf with pdf  $f_\theta(x) = f(x - \theta)$ , where  $f(x) = (1/2)e^{-|x|}$ ,  $\theta$  is unknown and we want to estimate  $\theta$ . Let  $T(F_\theta)$  be a solution to the equation

$$\lambda_{F_\theta}(t) = \int \Psi(x - t) dF_\theta(x) = 0 \text{ where } \Psi(x) = \begin{cases} |x|^{1/2} & x \leq 0 \\ -|x|^{1/2} & x > 1. \end{cases}$$

Then the influence function of  $T(F_\theta)$  is  $T'(x; F_\theta) = -\Psi(x - T(F_\theta)) / \lambda'_{F_\theta}(T(F_\theta))$ .

- (a) Show that  $T(F_\theta) = \theta$ .
  - (b) Express the corresponding  $M$ -estimator  $\hat{\theta}_n = T(F_n)$  in terms of  $(X_1, \dots, X_n)$  of which  $F_n$  is the edf.
  - (c) Write down one-term Taylor expansion of  $\sqrt{n}(\hat{\theta}_n - \theta)$ . Assume that the remainder term is  $o_P(1)$ .
  - (d) Show that  $E_\theta[T'(X; F_\theta)] = 0$  and evaluate  $\sigma^2(F_\theta, \Psi) = \text{Var}[T'(X; F_\theta)]$ .
  - (e) Find the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_n - \theta)$ . [Hint: To find  $\lambda'_{F_\theta}(T(F_\theta)) = d\lambda_{F_\theta}(t)/dt|_{t=\theta}$ , differentiate under the integral and use  $\Psi'(x) = -(1/2)|x|^{-1/2}$  for  $x \neq 0$ . Although  $\Psi'$  does not exist at  $x = 0$ , this calculation is still valid. Also note that  $\int_0^\infty x^{-1/2} e^{-x} dx = \Gamma(1/2) = \sqrt{\pi}$ .]
- 10.18.** Let  $f_\theta(x) = f(x - \theta)$  be the pdf of a double exponential distribution with mean  $\theta$  (ie,  $f_\theta(x) = (1/2)e^{-|x-\theta|}$ ) having cdf  $F_\theta$  and let  $F_n$  the empirical cdf of a random sample from  $F_\theta$ . Let  $T(F_n)$  be the  $L$ -estimator of  $\theta$  with score function  $J(u) = 4uI(0 \leq u \leq 1/2) + 4(1-u)I(1/2 < u \leq 1)$ .
- (a) Express  $T(F_n)$  as a linear function of the order statistics  $X_{n:1} < \dots < X_{n:n}$ , using coefficients  $J(i/(n+1))$  instead of  $J(i/n)$ , with  $n$  even.
  - (b) Show that  $T(F_\theta) = \theta$  for the corresponding  $L$ -functional.
  - (c) Find the asymptotic distribution of  $\sqrt{n}[T(F_n) - \theta]$ . [Hint: The asymptotic variance  $\sigma^2(F, J)$  of  $\sqrt{n}[T(F_n) - T(F)]$  involves  $F^{-1}$ . If  $F$  is the cdf corresponding to  $f(x) = (1/2)e^{-|x|}$ , then  $F^{-1}(u) = \log(2u)$  if  $0 \leq u \leq 1/2$  and  $= -\log(2(1-u))$  if  $1/2 < u \leq 1$ .]

# Linear Models

## 11.1 Introduction

Linear models are widely used in statistical data analysis when the dependent or the response variable is quantitative, whereas the independent variables may be quantitative, qualitative, or both. It can also be used for some types of nonlinear modeling as an example given below will show. A few well-known classes of linear models are

- (i) regression: all the variables are quantitative,
- (ii) analysis of variance (ANOVA): all the independent variables are qualitative, and
- (iii) analysis of covariance (ANCOVA): some of the independent variables are quantitative and some qualitative.

An obvious example of a linear model is simple linear regression with one independent variable. If the observations are  $(Y_i, X_i)$ ,  $i = 1, \dots, n$ , where  $\{Y_i\}$  are the values of the dependent variable and  $\{X_i\}$  are the values of the independent variable, then the simple linear regression model is

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\beta_0$  is the intercept,  $\beta_1$  is the slope, and  $\{\varepsilon_i\}$  are mutually uncorrelated random errors with mean 0 and common variance  $\sigma^2$ . This model may also be written as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix},$$

where  $E[\boldsymbol{\varepsilon}] = \mathbf{0}$  and  $\text{Cov}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}$ . In the linear model terminology,  $\mathbf{X}$  is called the design matrix and the goal is to obtain estimates of the unknown parameters  $\boldsymbol{\beta}$  and  $\sigma^2$ , and carry out inferences on them.

It turns out that all the models mentioned in (i)–(iii) can be rewritten in the framework of a Gauss-Markov model which is

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad E[\boldsymbol{\varepsilon}] = \mathbf{0}, \quad \text{Cov}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}, \quad (1)$$

where  $\mathbf{Y}$  is  $n \times 1$  vector of observed response values, the design matrix  $\mathbf{X}$  is of order  $n \times p$  and  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unknown parameters. In standard applications, the errors are often taken to be iid  $N(0, \sigma^2)$ . In this chapter, the columns of  $\mathbf{X}$  will sometimes be referred to as independent variables.

### Important Assumptions

Throughout this chapter, we assume that the design matrix  $\mathbf{X}$  is nonrandom, or if it is random all the calculations such as  $E[\cdot]$ ,  $\text{Cov}[\cdot]$ , etc., are carried out conditionally on  $\mathbf{X}$ . We also assume that  $n > p$  and  $\mathbf{X}$  has full rank (ie,  $\text{rank}(\mathbf{X}) = p$ ). For the Gauss-Markov model,

$$E[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta} \text{ and } \text{Cov}[\mathbf{Y}] = \sigma^2 \mathbf{I}.$$

The problems of inference involve estimation of  $\boldsymbol{\beta}$  and  $\sigma^2$ , construction of confidence intervals, and hypotheses tests for  $\boldsymbol{\beta}$  and its linear functions, deciding if some columns of  $\mathbf{X}$  can be deleted from the model, prediction of  $\mathbf{Y}$  at a future set of values of the independent variables, etc.

## 11.2 Examples of Gauss-Markov Models

Even though the description of the Gauss-Markov in Eq. (1) requires only the mutual uncorrelatedness of the random errors, this assumption is too general to be useful in applications. Therefore, following the standard practice, we assume that these random errors are iid with mean 0 and variance  $\sigma^2$  in all the examples below.

**Example 11.2.1** (Linear Regression). Suppose that we have  $n$  observation vectors  $(Y_i, X_{i,1}, \dots, X_{i,p-1})$ ,  $i = 1, \dots, n$ , where the response for the  $i$ th case is  $Y_i$  and values of the independent variables are  $X_{i,1}, \dots, X_{i,p-1}$ . Then a linear regression model is

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \dots + \beta_{p-1} X_{i,p-1} + \varepsilon_i,$$

where  $\{\varepsilon_i\}$  are iid with mean 0 and variance  $\sigma^2$ . The statistical analysis involves estimation of the unknown constants  $\beta_0, \beta_1, \dots, \beta_{p-1}$  and  $\sigma^2$  from the data  $(Y_i, X_{i,1}, \dots, X_{i,p-1})$ ,  $i = 1, \dots, n$ . This model can be expressed in the Gauss-Markov framework given in Eq. (1) with

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & X_{1,1} & \cdots & \cdots & \cdots & X_{1,p-1} \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 1 & X_{n,1} & \cdots & \cdots & \cdots & X_{n,p-1} \end{pmatrix},$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

**Example 11.2.2** (Nonlinear Regression). The structure given in the last example is quite general as it can accommodate nonlinear cases. For instance, if we have only one independent variable but it is believed that the relation between the independent variable and the dependent variable is nonlinear, then we may fit a polynomial model to account for the nonlinearity. If a polynomial of degree  $p - 1$  is considered, when the observations are  $\{Y_i, X_i\}$ ,  $i = 1, \dots, n$ , then we may consider the model

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \dots + \beta_{p-1} X_i^{p-1} + \varepsilon_i.$$

This is clearly of the form given in the last example if we take  $X_{i,1} = X_i, X_{i,2} = X_i^2, \dots, X_{i,p-1} = X_i^{p-1}$ . It is worthwhile to point out that in actual data analysis one may not use the powers of  $X_i$  to create independent variables since it may lead to the problem of very strong correlation among independent variables (also called the problem of high multicollinearity), which leads to instability (high variance) in the parameter estimates, especially when  $p$  is not small. In such cases, one may employ orthogonal polynomials to create the independent variables instead of using the powers of the independent variable to form the columns of the design matrix  $\mathbf{X}$ .

*Remark 11.2.1.* The last two examples show that any linear regression model can be reexpressed in the Gauss-Markov framework. In order to show that the ANOVA and ANCOVA models are in the Gauss-Markov setup, it will be enough to show that they can be written as linear regression models and this is the approach taken in the examples below.

**Example 11.2.3** (One-Factor Analysis of Variance (ANOVA)). The superintendent of a school district may be interested in comparing the mathematical aptitudes of the students in  $k$  different schools in a city. In order to achieve this,  $n_i$  students are selected at random from the  $i$ th school,  $i = 1, \dots, k$ , and the score of each of the  $n = n_1 + \dots + n_k$  students on a standardized mathematics test is recorded. This is an example of a one-factor study where “school” is called a factor with  $k$  levels. Thus there are  $k$  populations and  $n_i$  iid observations are available from the  $i$ th population. A typical assumption is that the populations may have different means  $\{\mu_i\}$  but the variances are the same. If  $Y_{ij}$  is the score of the  $j$ th student in the  $i$ th school, then the one-factor ANOVA model can be written as

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\{\varepsilon_{ij}\}$  are iid with mean 0 and variance  $\sigma^2$ . This model is called *balanced* if  $n_i$ 's are the same (ie,  $n_i = n_0$  for all  $i$ ).

If  $\mu$  is an overall (weighted) average of  $\{\mu_i\}$  (ie,  $\mu = \sum_{i=1}^k w_i \mu_i$  where  $w_i \geq 0$  and  $\sum w_i = 1$ ), then the factor effect for the  $i$ th school is defined to be  $\alpha_i = \mu_i - \mu$  and  $\{\alpha_i\}$  satisfy the constraint  $\sum w_i \alpha_i = 0$ . In practice, the weights are often taken to be  $w_i \equiv 1/k$  or  $w_i = n_i/n$ ,  $i = 1, \dots, k$ , though other choices are also possible. The one-factor ANOVA model may also be written as a factor-effects model

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k.$$

Note that  $E[Y_{ij}] = \mu + \alpha_i = \mu_i$ , irrespective of how  $\mu$  and  $\alpha_i$  are defined.

We now examine how this model may be recast in the regression setup. If the main effects satisfy the constraint  $\sum \alpha_i = 0$  (ie,  $w_i \equiv 1/k$ ), and if  $\alpha_1, \dots, \alpha_{k-1}$  are known, then  $\alpha_k = -\alpha_1 - \dots - \alpha_{k-1}$  is known. Thus there are really  $k$  free parameters in the factor-effects model and they are  $\mu, \alpha_1, \dots, \alpha_{k-1}$ . Define  $k-1$  variables as follows

$$X_{ij,1} = \begin{cases} 1 & j=1 \\ -1 & j=k \\ 0 & \text{otherwise} \end{cases}, \quad X_{ij,2} = \begin{cases} 1 & j=2 \\ -1 & j=k \\ 0 & \text{otherwise} \end{cases}, \dots,$$

$$X_{ij,k-1} = \begin{cases} 1 & j=k-1 \\ -1 & j=k \\ 0 & \text{otherwise.} \end{cases}$$

Then the ANOVA model can be expressed in the regression framework

$$Y_{ij} = \beta_0 + \beta_1 X_{ij,1} + \dots + \beta_{k-1} X_{ij,k-1} + \varepsilon_{ij}, \text{ with}$$

$$\beta_0 = \mu, \quad \beta_1 = \alpha_1, \dots, \beta_{k-1} = \alpha_{k-1}.$$

Note that when  $j = k$ ,

$$\mu + \alpha_1 X_{ij,1} + \dots + \alpha_{k-1} X_{ij,k-1} = \mu - (\alpha_1 + \dots + \alpha_{k-1}) = \mu + \alpha_k.$$

**Example 11.2.4** (One-Factor ANOVA Continued). As mentioned in the last example, there are many ways to define the overall mean  $\mu$  and the factor effects  $\{\alpha_i\}$ . If the overall mean  $\mu$  is defined to be  $\mu = \sum w_i \mu_i$  with  $w_i = n_i/n$  and  $\alpha_i = \mu_i - \mu$ , then the constraint on the factor effects is  $\sum (n_i/n) \alpha_i = 0$  (ie,  $\sum n_i \alpha_i = 0$ ). In this case one may define the  $X$ -variables as

$$X_{ij,1} = \begin{cases} 1 & i=1 \\ -n_1/n_k & i=k \\ 0 & \text{otherwise} \end{cases}, \quad X_{ij,2} = \begin{cases} 1 & i=2 \\ -n_2/n_k & i=k \\ 0 & \text{otherwise} \end{cases}, \dots,$$

$$X_{ij,k-1} = \begin{cases} 1 & i=k-1 \\ -n_{k-1}/n_k & i=k \\ 0 & \text{otherwise.} \end{cases}$$

Then the ANOVA model can be rewritten as

$$Y_{ij} = \beta_0 + \beta_1 X_{ij,1} + \dots + \beta_{k-1} X_{ij,k-1} + \varepsilon_{ij}, \text{ with}$$

$$\beta_0 = \mu, \quad \beta_1 = \alpha_1, \dots, \beta_{k-1} = \alpha_{k-1}.$$

A one-factor ANOVA model may also be written as a regression of  $Y$  on the following  $k-1$  indicator variables as defined below

$$X_{ij,1} = \begin{cases} 1 & i=1 \\ 0 & i \neq 1 \end{cases}, \dots, X_{ij,k-1} = \begin{cases} 1 & i=k-1 \\ 0 & i \neq k-1 \end{cases}.$$

**Example 11.2.5** (Two-Factor ANOVA). If in [Example 11.2.3](#), ethnicity/race of each student is recorded, then we have a two-factor study with factors “school” and “ethnicity.” Changing the notations a bit, suppose that there  $a$  levels of factor  $A$  (school) and  $b$  levels of factor  $B$  (ethnic group). Assume that a random sample of  $n_{ij}$  students is taken from the  $i$ th school with the  $j$ th ethnic background and the observed scores are  $\{Y_{ijk}: k=1, \dots, n_{ij}\}$ . If

$\mu_{ij}$  is the mean score of the students in the  $i$ th school with the  $j$ th ethnic background, then the cell means model is

$$Y_{ijk} = \mu_{ij} + \varepsilon_{ijk}, \quad k = 1, \dots, n_{ij}, \quad j = 1, \dots, b, \quad i = 1, \dots, a,$$

where  $\{\varepsilon_{ijk}\}$  are usually assumed to be iid random errors with mean 0 and variance  $\sigma^2$ . A two-factor study is called balanced if  $n_{ij} = n_0$  for all  $i$  and  $j$ .

In order to rewrite the two-factor model as a factor-effects model, let us define

$$\mu = (ab)^{-1} \sum_{j=1}^a \sum_{i=1}^b \mu_{ij}, \quad \mu_{i\cdot} = b^{-1} \sum_{j=1}^b \mu_{ij}, \quad \mu_{\cdot j} = a^{-1} \sum_{i=1}^a \mu_{ij},$$

$$\alpha_i = \mu_{i\cdot} - \mu, \quad \beta_j = \mu_{\cdot j} - \mu, \text{ and}$$

$$(\alpha\beta)_{ij} = \mu_{ij} - \mu_{i\cdot} - \mu_{\cdot j} + \mu = \mu_{ij} - (\mu + \alpha_i + \beta_j), \text{ and hence}$$

$$\mu_{ij} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij}.$$

Here  $\mu$  is the overall mean,  $\{\alpha_i\}$  are called the main effects of factor  $A$ ,  $\{\beta_j\}$  the main effects of factor  $B$ , and  $\{(\alpha\beta)_{ij}\}$  the interaction effects. The main effects and the interactions satisfy the constraints

$$\begin{aligned} \sum \alpha_i &= 0, & \sum \beta_j &= 0, \\ \sum_j (\alpha\beta)_{ij} &= 0 \text{ for any } i, & \text{and} \quad \sum_i (\alpha\beta)_{ij} &= 0 \text{ for any } j. \end{aligned}$$

Thus the two-factor ANOVA model may be written as a factor-effects model

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}.$$

If it turns out that the interaction effects are zero, then we end up with an additive model (ie, additive in factor effects)

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}.$$

In order to express the two-factor ANOVA model in the regression framework, we need to define variables for factors  $A$  and  $B$ . We can create  $a - 1$  factor  $A$  variables as

$$\begin{aligned} X_{ijk,1}^{(A)} &= \begin{cases} 1 & i = 1 \\ -1 & i = a \\ 0 & \text{otherwise} \end{cases}, & X_{ijk,2}^{(A)} &= \begin{cases} 1 & i = 2 \\ -1 & i = a \\ 0 & \text{otherwise} \end{cases}, \dots, \\ X_{ijk,a-1}^{(A)} &= \begin{cases} 1 & i = a - 1 \\ -1 & i = a \\ 0 & \text{otherwise} \end{cases}. \end{aligned}$$

Similarly, the  $b - 1$  factor  $B$  variables are

$$\begin{aligned} X_{ijk,1}^{(B)} &= \begin{cases} 1 & j = 1 \\ -1 & j = b \\ 0 & \text{otherwise} \end{cases}, & X_{ijk,2}^{(B)} &= \begin{cases} 1 & j = 2 \\ -1 & j = b \\ 0 & \text{otherwise} \end{cases}, \dots, \end{aligned}$$

$$X_{ijk,b-1}^{(B)} = \begin{cases} 1 & j = b-1 \\ -1 & j = b \\ 0 & \text{otherwise.} \end{cases}$$

Then the two-factor ANOVA model can be expressed as

$$\begin{aligned} Y_{ijk} = & \mu + \sum_{l=1}^{a-1} \alpha_l X_{ijk,l}^{(A)} + \sum_{m=1}^{b-1} \beta_m X_{ijk,m}^{(B)} \\ & + \sum_{l=1}^{a-1} \sum_{m=1}^{b-1} (\alpha\beta)_{lm} X_{ijk,l}^{(A)} X_{ijk,m}^{(B)} + \varepsilon_{ijk}. \end{aligned}$$

It is of interest to note that the interaction effects are the coefficients associated with the product terms of  $X^{(A)}$  and  $X^{(B)}$ .

The additive model (ie, when the interactions are not present) can be written as

$$Y_{ijk} = \mu + \sum_{l=1}^{a-1} \alpha_l X_{ijk,l}^{(A)} + \sum_{m=1}^{b-1} \beta_m X_{ijk,m}^{(B)} + \varepsilon_{ijk}.$$

**Example 11.2.6** (Analysis of Covariance). Analysis of covariance models come up when some of the independent variables are quantitative and others are qualitative. Let us first discuss a case with one qualitative variable (factor) and one quantitative variable. If in [Example 11.2.3](#), family income level of each student in the sample is recorded, then any modeling should take into account the school effect (factor) and the income levels  $\{Z_{ij}\}$ . A simple model in this case is

$$\begin{aligned} Y_{ij} &= \mu + \alpha_i + \gamma Z_{ij} + \varepsilon_{ij}, \text{ or more generally} \\ Y_{ij} &= \mu + \alpha_i + \gamma_i Z_{ij} + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k, \end{aligned}$$

where  $\{\varepsilon_{ij}\}$  are iid with mean 0 and variance  $\sigma^2$ . The first is an additive model, whereas the second model contains an interaction between the qualitative and the quantitative variables. In the first model, there are  $k$  parallel regression lines, whereas the second model allows for  $k$  separate regression lines with possibly different intercepts and different slopes. In both cases, it is assumed that  $\sum \alpha_i = 0$ . In order to write these two models in the Gauss-Markov framework, we may define the indicator variables  $X_{ij,1}, \dots, X_{ij,k}$  for  $k$  schools as in [Example 11.2.4](#), that is,

$$X_{ij,1} = \begin{cases} 1 & \text{if } i = 1 \\ 0 & \text{otherwise} \end{cases}, \dots, X_{ij,k} = \begin{cases} 1 & \text{if } i = k \\ 0 & \text{otherwise.} \end{cases}$$

Then the first model can be written as

$$Y_{ij} = \sum_{i=1}^k \beta_i X_{ij,i} + \gamma Z_{ij} + \varepsilon_{ij},$$

where  $\beta_l = \mu + \alpha_l$ . The second model (ie, the model with interactions) can be written as

$$Y_{ij} = \sum_{l=1}^k \beta_l X_{ij,l} + \sum_{l=1}^k \gamma_l X_{ij,l} Z_{ij} + \varepsilon_{ij},$$

with  $\beta_l = \mu + \alpha_l$ .

**Example 11.2.7** (Analysis of Covariance). A researcher wishes to investigate the effects of  $k$  different diets on the growth (weight) of animals. Let  $Y_{ij}$  be the growth of the  $j$ th subject on the  $i$ th diet and let  $Z_{ij1}$  and  $Z_{ij2}$  be the initial weight and age of the subject. Analysis of covariance model may be written as

$$Y_{ij} = \mu + \alpha_i + \gamma_1 Z_{ij1} + \gamma_2 Z_{ij2} + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\{\varepsilon_{ij}\}$  are iid with mean 0 and variance  $\sigma^2$ . If  $\{X_{ij,l}: l = 1, \dots, k\}$  are created as in the last example, then we may write this model as

$$Y_{ij} = \sum_{l=1}^k \beta_l X_{ij,l} + \sum_{l=1}^2 \gamma_l Z_{ijl} + \varepsilon_{ij}.$$

If the researcher wishes to consider, in addition to the diet (factor  $A$ ,  $a$  levels), the effect of gender (factor  $B$ ), then she may have  $n_{ij}$  subjects of gender  $j$  assigned to diet  $i$ . If  $Y_{ijk}$  is the growth rate of the  $k$ th subject with gender  $j$  and diet  $i$ , we may consider the model

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \gamma_1 z_{ijk1} + \gamma_2 Z_{ijk2} + \varepsilon_{ijk}, \quad k = 1, \dots, n_{ij}, \quad j = 1, 2, \quad i = 1, k,$$

where  $\{\alpha_i\}$ ,  $\{\beta_j\}$  are the main effects of the factors and  $\{(\alpha\beta)_{ij}\}$  the interaction effects, and they satisfy the constraints stated in [Example 11.2.5](#).

*Remark 11.2.2.* In the last example, the first model with diet as the factor, and initial weight and age as the covariates, may be written as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon},$$

where  $\mathbf{Z}$  is an  $n \times 2$  matrix whose columns consist of  $\{Z_{ij1}\}$  and  $\{Z_{ij2}\}$ . The first component on the right-hand side is  $\mathbf{X}\boldsymbol{\beta}$  and this consists of information on the factor levels. The second component  $\mathbf{Z}\boldsymbol{\gamma}$  consists of information on the covariates. Note that this follows the Gauss-Markov framework with the  $n \times (k+2)$  design matrix  $[\mathbf{X} \ \mathbf{Z}]$  and the unknown vector of parameters  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k, \gamma_1, \gamma_2)^T$ . The same representation holds for the model with interaction in [Example 11.2.6](#) (which allows for regression lines with different slopes and intercepts), except that  $\mathbf{Z}$  is now an  $n \times k$  matrix whose  $l$ th column has  $\{X_{ij,l} Z_{ij}\}$ ,  $l = 1, \dots, k$ .

### 11.3 Gauss-Markov Models: Estimation

This section is devoted to estimation of the unknown parameters of a Gauss-Markov model as given in Eq. (1) (ie, to estimate  $\boldsymbol{\beta}$  and  $\sigma^2$  when  $\mathbf{Y}$  and  $\mathbf{X}$  are observed). The method of least squares is a standard procedure for obtaining an estimate of  $\boldsymbol{\beta}$  and it is done by

minimizing the quantity  $G(\mathbf{b}) = \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|^2$  with respect to  $\mathbf{b}$  in  $\mathbb{R}^p$ . It turns out that  $G$  has a unique minimum and if the minimum is attained at  $\mathbf{b} = \hat{\beta}$ , then  $\hat{\beta}$  is taken to be an estimate of  $\beta$  and an estimate of the unknown mean vector  $\mathbf{X}\beta$  is  $\hat{Y} = \mathbf{X}\hat{\beta}$ . The vector of residuals  $\hat{\epsilon} = \mathbf{Y} - \hat{Y}$ , which is an estimate of the error vector  $\epsilon$ , can be used to estimate  $\sigma^2$ .

The column space of  $\mathbf{X}$  is  $\mathcal{M}(\mathbf{X}) = \{\mathbf{X}\mathbf{b}: \mathbf{b} \in \mathbb{R}^p\}$  and the (orthogonal) projection on it is given by  $\mathbf{Q}_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  (Section B.6). Note that  $\mathbf{Q}_X$  is symmetric and idempotent (ie,  $\mathbf{Q}_X^2 = \mathbf{Q}_X$ ), and  $\mathbf{I} - \mathbf{Q}_X$  is the projection on  $\mathcal{M}(\mathbf{X})^\perp$ , the orthogonal complement of  $\mathcal{M}(\mathbf{X})$ .

### 11.3.1 Estimation of $\beta$ and $\sigma^2$

We begin with the discussion on estimation of  $\beta$ . Note that

$$G(\mathbf{b}) = \mathbf{Y}^T \mathbf{Y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{Y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}.$$

The gradient and Hessian of  $G$  are

$$\partial G / \partial \mathbf{b} = -2\mathbf{X}^T \mathbf{Y} + 2\mathbf{X}^T \mathbf{X} \mathbf{b}, \quad \partial^2 G / \partial \mathbf{b} \partial \mathbf{b} = 2\mathbf{X}^T \mathbf{X}.$$

If  $\hat{\beta}$  is a solution of  $\partial G / \partial \mathbf{b} = \mathbf{0}$ , then clearly  $\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{Y}$ . Since  $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T \mathbf{X}) = p$ , the Hessian  $\partial^2 G / (\partial \mathbf{b} \partial \mathbf{b}) = 2\mathbf{X}^T \mathbf{X}$  is positive definite and the function  $G$  is strictly convex. Therefore,  $G$  has a unique minimum at  $\mathbf{b} = \hat{\beta}$ .

The estimated mean vector and the vector of the residuals are

$$\begin{aligned} \hat{Y} &= \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{Q}_X \mathbf{Y}, \text{ and} \\ \hat{\epsilon} &= \mathbf{Y} - \hat{Y} = \mathbf{Y} - \mathbf{X}\hat{\beta} = (\mathbf{I} - \mathbf{Q}_X)\mathbf{Y}. \end{aligned}$$

Since  $\hat{Y} = \mathbf{X}\hat{\beta}$  is in  $\mathcal{M}(\mathbf{X})$  and  $\hat{\epsilon}$  is in  $\mathcal{M}(\mathbf{X})^\perp$ ,  $\hat{Y}$  is orthogonal to the vector of residuals. Thus we are led to the following important result on the least squares method for estimating  $\beta$ .

**Theorem 11.3.1.** Consider the function  $G(\mathbf{b}) = \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|^2$ ,  $\mathbf{b} \in \mathbb{R}^p$ .

- (a) The function  $G$  has a unique minimum and denote by  $\hat{\beta}$  the vector at which  $G$  achieves its minimum. Then we have  $\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{Y}$  (ie,  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ ).
- (b) Let  $\hat{\epsilon} = \mathbf{Y} - \mathbf{X}\hat{\beta}$  be the vector of residuals. Then  $\mathbf{X}^T \hat{\epsilon} = 0$ .
- (c) The vector of residuals is orthogonal to the estimated mean vector  $\hat{Y} = \mathbf{X}\hat{\beta}$  (ie,  $\hat{Y}^T \hat{\epsilon} = 0$ ).

A few important properties of the least squares estimate  $\hat{\beta}$  of  $\beta$ , the estimated mean vector  $\hat{Y} = \mathbf{X}\hat{\beta}$ , and the residual vector  $\hat{\epsilon}$  can be derived rather easily using some basic algebra. Since  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , we have

$$\begin{aligned} \mathbb{E}[\hat{\beta}] &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\mathbf{Y}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta, \\ \text{Cov}[\hat{\beta}] &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Cov}[\mathbf{Y}] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}, \\ \mathbb{E}[\mathbf{X}\hat{\beta}] &= \mathbf{X}\mathbb{E}[\hat{\beta}] = \mathbf{X}\beta, \\ \text{Cov}[\mathbf{X}\hat{\beta}] &= \mathbf{X}\text{Cov}[\hat{\beta}]\mathbf{X}^T = \sigma^2 \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \sigma^2 \mathbf{Q}_X. \end{aligned}$$

Since  $\mathbf{I} - \mathbf{Q}_X$  is symmetric and idempotent, we have

$$\begin{aligned}\mathbb{E}[\hat{\boldsymbol{\varepsilon}}] &= (\mathbf{I} - \mathbf{Q}_X)\mathbb{E}[Y] = (\mathbf{I} - \mathbf{Q}_X)\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \\ \text{Cov}[\hat{\boldsymbol{\varepsilon}}] &= (\mathbf{I} - \mathbf{Q}_X)\text{Cov}[Y](\mathbf{I} - \mathbf{Q}_X)^T \\ &= \sigma^2(\mathbf{I} - \mathbf{Q}_X)(\mathbf{I} - \mathbf{Q}_X)^T = \sigma^2(\mathbf{I} - \mathbf{Q}_X)^2 = \sigma^2(\mathbf{I} - \mathbf{Q}_X), \\ \mathbb{E}[\|\hat{\boldsymbol{\varepsilon}}\|^2] &= \text{trace}(\text{Cov}[\hat{\boldsymbol{\varepsilon}}]) = \sigma^2 \text{trace}(\mathbf{I} - \mathbf{Q}_X) \\ &= \sigma^2[n - \text{trace}(\mathbf{Q}_X)] = \sigma^2(n - p),\end{aligned}$$

the last step is justified since the rank of the projection matrix  $\mathbf{Q}_X$  equals its trace. The last result indicates that an unbiased estimator of  $\sigma^2$  is given by  $\hat{\sigma}^2 = \|\hat{\boldsymbol{\varepsilon}}\|^2/(n - p)$ . In the literature,

- (i)  $\|\hat{\boldsymbol{\varepsilon}}\|^2$  is usually called the *residual sum of squares* and is denoted by *SSE*,
- (ii) degrees of freedom (df) of the SSE is defined to be  $n - \text{rank}(\mathbf{X}) = n - p$ ,
- (iii) mean square error (denoted by *MSE*) is defined to be  $\text{MSE} = \|\hat{\boldsymbol{\varepsilon}}\|^2/(n - p)$ .

For any linear function  $\mathbf{l}^T \mathbf{Y}$  of  $\mathbf{Y}$ ,  $\mathbf{l} \in \mathbb{R}^n$ , we have

$$\text{Cov}\left[\mathbf{l}^T \hat{\mathbf{Y}}, \hat{\boldsymbol{\varepsilon}}\right] = \text{Cov}\left[\mathbf{l}^T \mathbf{Q}_X \mathbf{Y}, (\mathbf{I} - \mathbf{Q}_X) \mathbf{Y}\right] = \sigma^2 \mathbf{l}^T \mathbf{Q}_X (\mathbf{I} - \mathbf{Q}_X) = \mathbf{0}.$$

Thus any linear function of the estimated mean vector  $\hat{\mathbf{Y}}$  is uncorrelated with the vector of residuals  $\hat{\boldsymbol{\varepsilon}}$ . This observation is crucial in inference since under the assumption of normality of  $\boldsymbol{\varepsilon}$ , uncorrelatedness implies independence and thus  $\mathbf{X}\hat{\boldsymbol{\beta}}$  is independent of  $\hat{\boldsymbol{\varepsilon}}$ . The discussion above leads to the following result.

**Theorem 11.3.2.** *Let  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  be the fitted mean vector and  $\hat{\boldsymbol{\varepsilon}}$  be the vector of residuals as in Theorem 11.3.1. Then the following hold:*

- (a)  $E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ ,  $\text{Cov}[\hat{\boldsymbol{\beta}}] = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ .
- (b)  $E[\hat{\mathbf{Y}}] = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{Cov}[\hat{\mathbf{Y}}] = \sigma^2 \mathbf{Q}_X$ .
- (c)  $E[\hat{\boldsymbol{\varepsilon}}] = \mathbf{0}$ ,  $\text{Cov}[\hat{\boldsymbol{\varepsilon}}] = \sigma^2(\mathbf{I} - \mathbf{Q}_X)$ .
- (d) The residual vector  $\hat{\boldsymbol{\varepsilon}}$  is uncorrelated to any linear function of the estimated mean vector  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ .
- (e)  $E[\text{MSE}] = \sigma^2$ , where  $\text{MSE} = \|\hat{\boldsymbol{\varepsilon}}\|^2/(n - p)$ .

### 11.3.2 Estimation of Linear Functions of $\boldsymbol{\beta}$

Often it is of interest to estimate linear functions of the unknown parameter  $\boldsymbol{\beta}$ . If  $\tilde{\boldsymbol{\beta}}$  is an unbiased for  $\boldsymbol{\beta}$ , then a linear function  $\mathbf{a}^T \tilde{\boldsymbol{\beta}}$ ,  $\mathbf{a} \in \mathbb{R}^p$ , is also unbiased for  $\mathbf{a}^T \boldsymbol{\beta}$ . Let  $\mathbf{L}$  be a known  $p \times m$  matrix of rank  $m \leq p$ , and consider the problem of estimating the linear function  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$ . Least squares estimate of  $\boldsymbol{\theta}$  is defined to be equal to  $\hat{\boldsymbol{\theta}} = \mathbf{L}^T \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}}$  is the least squares estimate of  $\boldsymbol{\beta}$ . The following simple result shows that  $\hat{\boldsymbol{\theta}}$  is an unbiased estimate of  $\boldsymbol{\theta}$ .

**Lemma 11.3.1.**  $E[\hat{\theta}] = \theta$ ,  $Cov[\hat{\theta}] = \sigma^2 \mathbf{L}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}$ .

### 11.3.3 Best Linear Unbiased Estimation

We begin with a definition.

**Definition 11.3.1.** A linear function  $\tilde{\beta}$  of  $\mathbf{Y}$  is called a best linear unbiased estimator (BLUE) of  $\beta$  if

- (i)  $\tilde{\beta}$  is an unbiased estimator of  $\beta$ , and
- (ii) for any  $\mathbf{a} \in \mathbb{R}^p$ ,  $\text{Var}[\mathbf{a}^T \tilde{\beta}] \leq \text{Var}[\mathbf{l}^T \mathbf{Y}]$  for all linear unbiased estimators  $\mathbf{l}^T \mathbf{Y}$  of  $\mathbf{a}^T \beta$ ,  $\mathbf{l} \in \mathbb{R}^n$ .

It is clear from this definition that  $\tilde{\beta}$  is a BLUE of  $\beta$  if  $\mathbf{a}^T \tilde{\beta}$  is BLUE of  $\mathbf{a}^T \beta$  for any  $\mathbf{a} \in \mathbb{R}^p$ . The following argument will show that the BLUE exists, it is unique and it is equal to the least squares estimate  $\hat{\beta}$ . It will be enough to show that, if for any  $\mathbf{a} \in \mathbb{R}^p$ ,  $\mathbf{l}^T \mathbf{Y}$  is a BLUE for  $\mathbf{a}^T \beta$ , then  $\mathbf{l}^T \mathbf{Y} = \mathbf{a}^T \hat{\beta}$ .

If  $\mathbf{l}^T \mathbf{Y}$  is an unbiased estimator of  $\mathbf{a}^T \beta$ , then  $\mathbf{a}^T \beta = E[\mathbf{l}^T \mathbf{Y}] = \mathbf{l}^T \mathbf{X} \beta$  for all  $\beta$  and hence  $\mathbf{X}^T \mathbf{l} = \mathbf{a}$ . If  $\mathbf{l}^T \mathbf{Y}$  is a BLUE of  $\mathbf{a}^T \beta$ , then for any linear unbiased estimator  $\mathbf{m}^T \mathbf{Y}$  of 0 (ie,  $E[\mathbf{m}^T \mathbf{Y}] = 0$  for all  $\beta$ ),  $(\mathbf{l} + t\mathbf{m})^T \mathbf{Y}$  is also unbiased for  $\mathbf{a}^T \beta$ , where  $t$  is a real number. Let

$$h(t) = \text{Var}[(\mathbf{l} + t\mathbf{m})^T \mathbf{Y}] = \sigma^2 \|\mathbf{l} + t\mathbf{m}\|^2.$$

Since  $\mathbf{l}^T \mathbf{Y}$  is a BLUE, the function  $h$  achieves a minimum at  $t = 0$ , thus  $0 = h'(0) = 2\sigma^2 \mathbf{l}^T \mathbf{m}$  (ie,  $\mathbf{l}^T \mathbf{m} = 0$ ). Since  $\mathbf{m}^T \mathbf{Y}$  is an unbiased estimator of 0, we have  $\mathbf{m}^T \mathbf{X} \beta = 0$  for all  $\beta$  and thus  $\mathbf{X}^T \mathbf{m} = \mathbf{0}$ . Since  $\mathbf{l}^T \mathbf{m} = 0$  for all  $\mathbf{m}$  satisfying the condition  $\mathbf{X}^T \mathbf{m} = \mathbf{0}$  (ie, for all  $\mathbf{m} \in \mathcal{M}(\mathbf{X})^T$ ), it follows that  $\mathbf{l}$  must be in  $\mathcal{M}(\mathbf{X})$ . Thus  $\mathbf{l} = \mathbf{Xc}$  for some  $\mathbf{c} \in \mathbb{R}^p$ . Since  $\mathbf{l}^T \mathbf{Y}$  is unbiased for  $\mathbf{a}^T \beta$ , we have  $\mathbf{a}^T \beta = \mathbf{l}^T \mathbf{X} \beta = \mathbf{c}^T \mathbf{X}^T \mathbf{X} \beta$  for all  $\beta$ . This implies that  $\mathbf{c} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}$  and hence  $\mathbf{l} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}$ . Thus if  $\mathbf{l}^T \mathbf{Y}$  is a BLUE of  $\mathbf{a}^T \beta$ , then

$$\mathbf{l}^T \mathbf{Y} = \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{a}^T \hat{\beta}.$$

Uniqueness is clear since any BLUE must have this form.

Thus we are led to the following important result.

**Lemma 11.3.2.** *The BLUE of  $\beta$  is unique and it is equal to the least squares estimate  $\hat{\beta}$ .*

### 11.4 Decomposition of Total Sum of Squares

In each of the examples given in [Section 11.2](#) of this chapter, the design matrix  $\mathbf{X}$  has a column consisting of 1's. In regression, it corresponds to the intercept term and in ANOVA models, it corresponds to the overall mean. Let us assume without loss of generality that the first column of  $\mathbf{X}$  consists of 1's and the vector of parameters is  $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})^T$ . A model which only has the first column (and ignores the last  $p-1$  columns) is  $Y_i = \beta_0 + \varepsilon_i$ ,  $i = 1, \dots, n$ . Clearly the least squares estimate of  $\beta_0$  is  $\bar{Y} = n^{-1} \sum_{i=1}^n Y_i$  and the residual

sum of squares of this simple model  $\sum(Y_i - \hat{\beta}_0)^2 = \sum(Y_i - \bar{Y})^2$  is usually called the total sum of squares (SSTO). Let  $\hat{Y} = X\hat{\beta}$  be the least squares estimate of the mean  $X\beta$  for the full model  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$ . We already know that the residual vector  $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - X\hat{\beta} = \mathbf{Y} - \hat{Y}$  is orthogonal to any vector in  $\mathcal{M}(X)$  and hence to  $\bar{Y}\mathbf{1}$ , where  $\mathbf{1}$  is an  $n \times 1$  vector of 1's. Then

$$\begin{aligned} SSTO &= \sum(Y_i - \bar{Y})^2 = \|\mathbf{Y} - \bar{Y}\mathbf{1}\|^2 = \|X\hat{\beta} - \bar{Y}\mathbf{1}\|^2 + \|\mathbf{Y} - X\hat{\beta}\|^2 \\ &= \sum(\hat{Y}_i - \bar{Y})^2 + \sum(Y_i - \hat{Y}_i)^2 \\ &= \sum(\hat{Y}_i - \bar{Y})^2 + SSE. \end{aligned}$$

Since SSTO is the residual sum of squares for the model  $Y_i = \beta_0 + \varepsilon_i$ ,  $i = 1, \dots, n$ , the quantity  $\sum(\hat{Y}_i - \bar{Y})^2$  is the reduction in the residual sum of squares when we go from the simple model  $\mathbf{Y} = \beta_0\mathbf{1} + \boldsymbol{\varepsilon}$  to the full model  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$  and

$$R^2 = \frac{\sum(\hat{Y}_i - \bar{Y})^2}{SSTO} = 1 - SSE/SSTO$$

is the proportional reduction in the residual sum of squares. It is also called the *coefficient of determination* and is widely used in practice. Clearly,  $R^2$  is between 0 and 1, and if  $R^2$  is close 1, it is taken as an indication that the full model  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$  explains the data well. The reduction  $\sum(\hat{Y}_i - \bar{Y})^2$  is called the “regression sum of squares” in the regression model and “treatment sum of squares” in the one-factor ANOVA model.

Even though  $R^2$  is popular as a descriptive measure it has some flaws. If there is a true model  $\mathbf{Y} = X_*\beta_* + \boldsymbol{\varepsilon}$  with  $\text{rank}(X_*) = p_*$  and it is nested in the model under consideration  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$ , where  $\mathcal{M}(X_*) \subset \mathcal{M}(X)$  and  $\text{rank}(X) = p > p_*$ , then clearly the SSE for the latter model is smaller than that of the former (ie, the true lower dimensional model), and thus the latter model has a higher  $R^2$ . As a matter of fact, for any class of nested models, the value of  $R^2$  will always increase as we consider higher dimensional models, and this is true regardless of what the true model is. In order to remedy this, we first need to identify a parameter  $\rho^2$  that  $R^2$  is trying to estimate. It turns out that a descriptive measure called the adjusted  $R^2$ , denoted by  $R_{\text{adj}}^2$ , is a better estimate of  $\rho^2$  than  $R^2$  is. Before we go any further, let us first state a simple lemma and then define the concept of degrees of freedom associated with residual sums of squares for any linear model.

We consider here a general structure for the  $n \times 1$  observation vector  $\mathbf{Y}$  with mean  $\mu$  such that

$$\mathbf{Y} = \mu + \boldsymbol{\varepsilon}, \text{ with } E[\boldsymbol{\varepsilon}] = \mathbf{0} \quad \text{and} \quad \text{Cov}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}.$$

If  $X$  is  $n \times p$  of rank  $p$ , then the minimum of  $\|\mu - X\beta\|^2$ ,  $\beta \in \mathbb{R}^p$ , is attained at  $\beta^* = (X^T X)^{-1} X^T \mu$ , and the *projected mean* is  $X\beta^* = X(X^T X)^{-1} X^T \mu = Q_X \mu$ . So when a model  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$  is fitted to the data, then  $Q_X \mathbf{Y}$  is estimating the projected mean  $Q_X \mu$  and the model is a true description of the data if  $\mu = Q_X \mu$ .

**Lemma 11.4.1.** *Assume that  $\mathbf{Y} = \mu + \boldsymbol{\varepsilon}$ , where  $\mathbf{Y}$  is  $n \times 1$  observation vector,  $E[\boldsymbol{\varepsilon}] = \mathbf{0}$ ,  $\text{Cov}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}$ , and a model of the form  $\mathbf{Y} = X\beta + \boldsymbol{\varepsilon}$ , where  $X$  is  $n \times p$  of rank  $p$ , is fitted to the*

data. A measure of deviation of  $\mathbf{Y}$  from the projected mean  $\mathbf{Q}_X\boldsymbol{\mu}$  is  $D = E[\|\mathbf{Y} - \mathbf{Q}_X\boldsymbol{\mu}\|^2]$  and an estimate of  $D$  is  $SSE = \|\mathbf{Y} - \mathbf{Q}_X\mathbf{Y}\|^2$ . The following hold:

- (a)  $D = n\sigma^2 + \|(\mathbf{I} - \mathbf{Q}_X)\boldsymbol{\mu}\|^2$ ,
- (b)  $E[SSE] = (n - p) + \|(\mathbf{I} - \mathbf{Q}_X)\boldsymbol{\mu}\|^2$ .

**Definition 11.4.1.** If a model of the form  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{X}$  is  $n \times p$  matrix of rank  $p$ , is fitted to the observation vector  $\mathbf{Y}$  where  $\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$ ,  $E[\boldsymbol{\epsilon}] = \mathbf{0}$ ,  $\text{Cov}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{I}$ , then  $\mathbf{Q}_X\boldsymbol{\mu}$  is called the projected mean of the model. The degrees of freedom (df) of the residual sum of squares  $SSE = \|\mathbf{Y} - \mathbf{Q}_X\mathbf{Y}\|^2$  is  $n - p$ , and the mean square error is  $MSE = SSE/(n - p)$ .

For the simple model  $Y_i = \beta_0 + \varepsilon_i$ ,  $i = 1, \dots, n$ , the residual sum of squares is  $SSTO$  and its mean square error is  $MSTO/(n - 1)$ . The adjusted  $R^2$  is defined as

$$R_{\text{adj}}^2 = 1 - MSE/MSTO.$$

For the model  $\mathbf{Y} = \beta_0\mathbf{1} + \boldsymbol{\epsilon}$ , let  $\mathbf{Q}_0\boldsymbol{\mu}$  be the projection of  $\boldsymbol{\mu}$  on  $\mathbf{1}$  (ie,  $\mathbf{Q}_0\boldsymbol{\mu} = \beta_0\mathbf{1}$  with  $\beta_0 = \mathbf{1}^T\boldsymbol{\mu}/n$ ). Noting that  $\bar{Y}$  estimates  $\beta_0$  and  $SSTO$  estimates  $E[\|\mathbf{Y} - \mathbf{Q}_0\boldsymbol{\mu}\|^2]$ , the deviation of  $\mathbf{Y}$  from  $\mathbf{Q}_0\mathbf{X}\boldsymbol{\beta} = \beta_0\mathbf{1}$ . Similarly,  $SSE$  is an estimate of  $E[\|\mathbf{Y} - \mathbf{Q}_X\boldsymbol{\mu}\|^2]$ , when the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  (where the first column of  $\mathbf{X}$  consists of 1's) is fitted to the data. We may thus define the proportional reduction in the true deviation of  $\mathbf{Y}$  from the projected mean when we go from the model  $\mathbf{Y} = \beta_0\mathbf{1} + \boldsymbol{\epsilon}$  to the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  as

$$\begin{aligned} \rho^2 &= 1 - E[\|\mathbf{Y} - \mathbf{Q}_X\boldsymbol{\mu}\|^2]/E[\|\mathbf{Y} - \mathbf{Q}_0\boldsymbol{\mu}\|^2] \\ &= 1 - \frac{n\sigma^2 + \|(\mathbf{I} - \mathbf{Q}_X)\boldsymbol{\mu}\|^2}{n\sigma^2 + \|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2} \\ &= \frac{\|(\mathbf{Q}_X - \mathbf{Q}_0)\boldsymbol{\mu}\|^2}{n\sigma^2 + \|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2}, \end{aligned}$$

using the results in the lemma above. Note that  $R^2$  is an estimate of  $\rho^2$ , but it has a serious flaw as discussed above.

Under the assumption of normality of the error terms,  $\|\mathbf{Y} - \mathbf{Q}_0\boldsymbol{\mu}\|^2/\sigma^2 \sim \chi_{n-1}^2(\delta_0^2)$ , where the noncentrality parameter is  $\delta_0^2 = (1/2)\|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2/\sigma^2$ . Similarly,  $\|\mathbf{Y} - \mathbf{Q}_X\mathbf{Y}\|^2/\sigma^2 \sim \chi_{n-p}^2(\delta^2)$  with  $\delta^2 = (1/2)\|(\mathbf{I} - \mathbf{Q}_X)\boldsymbol{\mu}\|^2/\sigma^2$ . If  $n$  is large, it can be shown that

$$\begin{aligned} MSE &= \|\mathbf{Y} - \mathbf{Q}_X\mathbf{Y}\|^2/(n - p) \\ &= \sigma^2 + n^{-1}\|(\mathbf{I} - \mathbf{Q}_X)\boldsymbol{\mu}\|^2 + O_P(n^{-1/2}), \text{ and} \\ MSTO &= \|\mathbf{Y} - \mathbf{Q}_0\mathbf{Y}\|^2/(n - 1) \\ &= \sigma^2 + n^{-1}\|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2 + O_P(n^{-1/2}), \end{aligned}$$

assuming that the quantity  $n^{-1}\|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2$  stays bounded as  $n \rightarrow \infty$ . These results are generally true under reasonable technical conditions even if the error terms are not normally distributed. Thus

$$\begin{aligned}
R_{\text{adj}}^2 &= 1 - \text{MSE}/\text{MSTO} = \frac{\text{MSTO} - \text{MSE}}{\text{MSTO}} \\
&= \frac{n^{-1} \|(\mathbf{Q}_X - \mathbf{Q}_0)\boldsymbol{\mu}\|^2}{\sigma^2 + n^{-1} \|(\mathbf{I} - \mathbf{Q}_0)\boldsymbol{\mu}\|^2} + O_p(n^{-1/2}) \\
&= \rho^2 + O_p(n^{-1/2}).
\end{aligned}$$

Clearly,  $R_{\text{adj}}^2$  is a  $\sqrt{n}$  consistent estimate of  $\rho^2$ .

**Example 11.4.1.** In the multiple regression case with  $p - 1$  independent variables, the normal equations  $\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{Y}$  are

$$\begin{bmatrix} n & \sum X_{i,1} & \cdots & \sum X_{i,p-1} \\ \sum X_{i,1} & \sum X_{i,1}^2 & \cdots & \sum X_{i,1} X_{i,p-1} \\ \vdots & \vdots & \cdots & \vdots \\ \sum X_{i,p-1} & \sum X_{i,p-1} X_{i,1} & \cdots & \sum X_{i,p-1}^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} = \begin{bmatrix} \sum X_{i,1} Y_i \\ \sum X_{i,2} Y_i \\ \vdots \\ \sum X_{i,p-1} Y_i \end{bmatrix},$$

where all the sums are over  $i$  from 1 through  $n$ . When  $p = 2$  (ie, there is only one independent variable), the solutions are

$$\hat{\beta}_1 = \sum (X_{i,1} - \bar{X}_1)(Y_i - \bar{Y}) / \sum (X_{i,1} - \bar{X}_1)^2, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}_1,$$

where  $\bar{X}_1 = n^{-1} \sum_{i=1}^n X_{i,1}$  and  $\bar{Y} = n^{-1} \sum_{i=1}^n Y_i$ . When  $p > 2$ , there are no such simple expressions for the estimates of  $\beta_0$ ,  $\beta_1$ , etc. Typically, solving these equations require computing packages, which are widely available. If  $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$  (ie,  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{i,1} + \cdots + \hat{\beta}_{p-1} X_{i,p-1}$ ,  $i = 1, \dots, n$ ), then the quantity  $\sum (\hat{Y}_i - \bar{Y})$  is called the “regression sum of squares” and thus  $R^2 = \text{SSR}/\text{SSTO}$ .

**Example 11.4.2.** In the one-factor ANOVA case as in [Example 11.2.3](#), it is fairly easy to get the least squares estimate of  $\mu_i$

$$\hat{\mu}_i = n_i^{-1} \sum_{j=1}^{n_i} Y_{ij} := \bar{Y}_i.$$

The residual sum of squares is  $\text{SSE} = \sum_i \sum_j (Y_{ij} - \bar{Y}_{i..})^2$  and  $\text{MSE} = \text{SSE}/(n - k)$  is an unbiased estimate of  $\sigma^2$ , where  $n = n_1 + \cdots + n_k$  is the total number of observations. The estimates of the overall mean  $\mu = \sum w_i \mu_i$  and  $\alpha_i = \mu_i - \mu$  are obtained by substituting  $\{\mu_i\}$  by  $\{\hat{\mu}_i\}$ . For instance, if  $w_i = n_i/n$ , then

$$\hat{\mu} = \sum_i (n_i/n) \bar{Y}_{i..} = \sum_i \sum_j Y_{ij}/n := \bar{Y}_{...} \quad \text{and} \quad \hat{\alpha}_i = \bar{Y}_{i..} - \bar{Y}_{...}.$$

As before if we denote  $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$ , then  $\hat{Y}_{ij} = \bar{Y}_{i..}$  and  $\sum_i \sum_j (\hat{Y}_{ij} - \bar{Y}_{i..})^2 = \sum n_i (\bar{Y}_{i..} - \bar{Y}_{...})^2$  is called the “treatment sum of squares” (SSTR) or the “between group sum of squares.” For this case, the SSE is sometimes called the “within group sum of squares.” The

decomposition of the total sum of squares is

$$\begin{aligned} SSTO &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2 = \sum n_i \hat{\alpha}_i^2 + SSE, \text{ ie,} \\ SSTO &= SSTR + SSE. \end{aligned}$$

In this case  $R^2 = SSTR/SSTO$ . The quantity  $SSTR/(k-1)$  is usually called *MSTR*, mean square for the treatment. It can be shown that (left as an exercise)

$$E[MSTR]/\sigma^2 = 1 + (k-1)^{-1} \sum n_i (\mu_i - \mu)^2 / \sigma^2.$$

The quantity  $(k-1)^{-1} \sum n_i (\mu_i - \mu)^2 / \sigma^2$  is a unit-free measure of the variability of  $\{\mu_i\}$ . This measure equals 0 if and only if  $\mu_i$ 's are all the same. Thus the ratio  $F = MSTR/MSE$ , which fluctuates about 1 if and only if  $\mu_1 = \dots = \mu_k$ , is used for testing the hypothesis that the means are the same. Under the assumption of normality (ie,  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ ),

$$\begin{aligned} SSE/\sigma^2 &\sim \chi_{n-k}^2, \quad SSTR/\sigma^2 \sim \chi_{k-1}^2(\delta^2), \text{ where} \\ \delta^2 &= (1/2) \sum n_i \hat{\alpha}_i^2 / \sigma^2, \end{aligned}$$

and  $F = MSTR/MSE \sim F_{k-1, n-k}(\delta^2)$ . Hence the *F*-statistic can be used to test  $H_0: \alpha_1 = \dots = \alpha_k = 0$ , since  $F \sim F_{k-1, n-k}$  under  $H_0$ .

If the overall mean  $\mu$  is defined as  $k^{-1} \sum \mu_i$ , then the estimates of  $\mu$  and  $\alpha_i$  are  $\hat{\mu} = k^{-1} \sum \bar{Y}_{ij}$  and  $\hat{\alpha}_i = \bar{Y}_{ij} - \hat{\mu}$ . Note that  $E[Y_{ij}]$  is always equal to  $\mu_i$  irrespective of how  $\mu$  and  $\{\alpha_i\}$  are defined,  $\hat{Y}_{ij} = \bar{Y}_{ij}$  and  $\hat{\varepsilon}_{ij} = Y_{ij} - \bar{Y}_{ij}$ , and hence the residual sum of squares also remains the same.

**Example 11.4.3.** In general ANOVA models (one- or multifactor), one may first obtain the estimates of the means of all the factor combinations and then use them to estimate the overall mean, factor effects, interactions, etc. For instance, in the two-factor ANOVA model, the estimate of  $\mu_{ij}$  is  $\hat{\mu}_{ij} = \bar{Y}_{ij}$ , where  $\bar{Y}_{ij} = n_{ij}^{-1} \sum_{k=1}^{n_{ij}} Y_{ijk}$ , and this can be used to estimate  $\mu$ ,  $\mu_i$ ,  $\mu_{.j}$ , the factor effects  $\{\alpha_i\}$ ,  $\{\beta_j\}$ , and the interactions  $\{(\alpha\beta)_{ij}\}$ , which are all linear functions of  $\{\mu_{ij}\}$ . The residual sum of squares is  $SSE = \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij})^2$ , and  $MSE = SSE/(n - ab)$ , where  $n = \sum_i \sum_j n_{ij}$  is the total number of observations, is an unbiased estimate of  $\sigma^2$ . Thus

$$\begin{aligned} \hat{\mu} &= (ab)^{-1} \sum_i \sum_j \hat{\mu}_{ij} = (ab)^{-1} \sum_i \sum_j \bar{Y}_{ij}, \\ \hat{\mu}_{i.} &= b^{-1} \sum_j \bar{Y}_{ij}, \quad \hat{\mu}_{.j} = a^{-1} \sum_i \bar{Y}_{ij}, \\ \hat{\alpha}_i &= \hat{\mu}_{i.} - \hat{\mu}, \quad \hat{\beta}_j = \hat{\mu}_{.j} - \hat{\mu}, \quad \widehat{(\alpha\beta)}_{ij} = \hat{\mu}_{ij} - \hat{\mu}_{i.} - \hat{\mu}_{.j} + \hat{\mu}, \text{ and} \\ \hat{Y}_{ijk} &= \bar{Y}_{ij} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \widehat{(\alpha\beta)}_{ij}. \end{aligned}$$

**Example 11.4.4.** In the balanced two-factor ANOVA model (ie,  $n_{ij} \equiv n_0$ ), the estimates are much simpler when one uses the following notations

$$\begin{aligned}\bar{Y}_{...} &= (n_0 ab)^{-1} \sum_i \sum_j \sum_k Y_{ijk}, \\ \bar{Y}_{i..} &= (n_0 b)^{-1} \sum_j \sum_k Y_{ijk}, \quad \bar{Y}_{.j.} = (n_0 a)^{-1} \sum_i \sum_k Y_{ijk}, \text{ and} \\ \bar{Y}_{ij.} &= n_0^{-1} \sum_k Y_{ijk}.\end{aligned}$$

With these notations

$$\begin{aligned}\hat{\mu} &= \bar{Y}_{...}, \quad \hat{\alpha}_i = \bar{Y}_{i..} - \bar{Y}_{...}, \quad \hat{\beta}_j = \bar{Y}_{.j.} - \bar{Y}_{...}, \\ (\widehat{\alpha\beta})_{ij} &= \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...}, \text{ and} \\ Y_{ijk} - \hat{\mu} &= \hat{\alpha}_i + \hat{\beta}_j + (\widehat{\alpha\beta})_{ij} + \hat{\varepsilon}_{ijk},\end{aligned}$$

where  $\{\hat{\varepsilon}_{ijk} = Y_{ijk} - \bar{Y}_{ij.}\}$  are the residuals. If both sides are squared and summed over  $i, j$ , and  $k$ , all the cross-product terms vanish to yield the following decomposition of the total sum of squares

$$\begin{aligned}SSTO &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{...})^2 \\ &= \sum_i \sum_j \sum_k (\bar{Y}_{ij.} - \bar{Y}_{...})^2 + \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij.})^2 \\ &= (n_0 b) \sum_i \hat{\alpha}_i^2 + (n_0 a) \sum_j \hat{\beta}_j^2 + n_0 \sum_i \sum_j (\widehat{\alpha\beta})_{ij}^2 + SSE \\ &:= SSA + SSB + SSAB + SSE.\end{aligned}$$

$SSA$ ,  $SSB$ , and  $SSAB$  are called the sums of squares due to the main effects of factor  $A$ , main effects of factor  $B$ , and the interactions, respectively. It is important to point out that this decomposition of  $SSTO$  is no longer valid for the unbalanced case.

It can also be shown that (left as an exercise)

$$\begin{aligned}E[SSA] &= (a-1)\sigma^2 + (n_0 b) \sum_i \alpha_i^2, \quad E[SSB] = (b-1)\sigma^2 + (n_0 a) \sum_j \beta_j^2, \\ E[SSAB] &= (a-1)(b-1)\sigma^2 + n_0 \sum_i \sum_j (\widehat{\alpha\beta})_{ij}^2, \\ E[SSE] &= (n-ab)\sigma^2.\end{aligned}$$

How do we view these results in matrix terms? Consider the following (sub)models

$$\begin{aligned}Y_{ijk} &= \mu + \varepsilon_{ijk} \text{ (model 0),} \\ Y_{ijk} &= \mu + \alpha_i + \varepsilon_{ijk} \text{ (model 1),} \\ Y_{ijk} &= \mu + \beta_j + \varepsilon_{ijk} \text{ (model 2),} \\ Y_{ijk} &= \mu + \alpha_i + \beta_j + \varepsilon_{ijk} \text{ (model 3), and} \\ Y_{ijk} &= \mu + \alpha_i + \beta_j + (\widehat{\alpha\beta})_{ij} + \varepsilon_{ijk} \text{ (model 4),}\end{aligned}$$

where model 4 is the full (true) model and the rest are submodels of the full model. The fitted values (or estimated means) of these models are

$$\begin{aligned}\hat{Y}_{ijk}^{(0)} &= \hat{\mu} = \bar{Y}_{...}, \quad \hat{Y}_{ijk}^{(1)} = \bar{Y}_{i..} = \hat{\mu} + \hat{\alpha}_i, \\ \hat{Y}_{ijk}^{(2)} &= \bar{Y}_{j..} = \hat{\mu} + \hat{\beta}_j, \quad \hat{Y}_{ijk}^{(3)} = \bar{Y}_{i..} + \bar{Y}_{j..} - \bar{Y}_{...} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j, \text{ and} \\ \hat{Y}_{ijk}^{(4)} &= \bar{Y}_{ij..} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \widehat{(\alpha\beta)}_{ij}.\end{aligned}$$

For each of the five linear models above, the postulated mean belongs to a linear space and the vector of fitted values is a projection of  $\mathbf{Y}$  on this column space. If  $\mathbf{Q}_j\mathbf{Y}$  denotes the projection of  $\mathbf{Y}$  for model  $j$ ,  $j = 0, \dots, 4$ , then the decomposition  $Y_{ijk} - \hat{\mu} = \hat{\alpha}_i + \hat{\beta}_j + \widehat{(\alpha\beta)}_{ij} + \hat{\varepsilon}_{ijk}$  can be written in the matrix form as

$$\begin{aligned}(\mathbf{I} - \mathbf{Q}_0)\mathbf{Y} &= (\mathbf{Q}_1 - \mathbf{Q}_0)\mathbf{Y} + (\mathbf{Q}_2 - \mathbf{Q}_0)\mathbf{Y} + (\mathbf{Q}_4 - \mathbf{Q}_3)\mathbf{Y} + (\mathbf{I} - \mathbf{Q}_4)\mathbf{Y}, \text{ or} \\ (\mathbf{I} - \mathbf{Q}_0)\mathbf{Y} &= \mathbf{M}_1\mathbf{Y} + \mathbf{M}_2\mathbf{Y} + \mathbf{M}_3\mathbf{Y} + \mathbf{M}_4\mathbf{Y},\end{aligned}$$

where  $\mathbf{M}_1 = \mathbf{Q}_1 - \mathbf{Q}_0$ ,  $\mathbf{M}_2 = \mathbf{Q}_2 - \mathbf{Q}_0$ ,  $\mathbf{M}_3 = \mathbf{Q}_4 - \mathbf{Q}_3$ , and  $\mathbf{M}_4 = \mathbf{I} - \mathbf{Q}_4$ . It is fairly easy to check that  $\mathbf{M}_j$ ,  $j = 1, 2, 3, 4$ , are projection matrices. Moreover, for the balanced two-factor ANOVA model,  $\mathbf{M}_i\mathbf{M}_j = 0$ ,  $1 \leq i \neq j \leq 4$ . When  $\{\varepsilon_{ijk}\}$  are iid  $N(0, \sigma^2)$ , using the results in Section B.7 we have

$$\begin{aligned}SSA/\sigma^2 &= \|\mathbf{M}_1\mathbf{Y}\|^2/\sigma^2 \sim \chi_{a-1}^2(\delta_1^2), \quad \delta_1^2 = (1/2)(n_0b) \sum \alpha_i^2/\sigma^2, \\ SSB/\sigma^2 &= \|\mathbf{M}_2\mathbf{Y}\|^2/\sigma^2 \sim \chi_{b-1}^2(\delta_2^2), \quad \delta_2^2 = (1/2)(n_0a) \sum \beta_j^2/\sigma^2, \\ SSAB/\sigma^2 &= \|\mathbf{M}_3\mathbf{Y}\|^2/\sigma^2 \sim \chi_{(a-1)(b-1)}^2(\delta_3^2), \quad \delta_3^2 = (1/2)n_0 \sum_i \sum_j (\alpha\beta)_{ij}^2/\sigma^2, \\ SSE/\sigma^2 &= \|\mathbf{M}_4\mathbf{Y}\|^2/\sigma^2 \sim \chi_{n-ab}^2,\end{aligned}$$

and SSA, SSB, SSAB, and SSE are independent.

**Example 11.4.5** (Estimation in One-Factor ANOVA). In one-factor ANOVA or multifactor models, one is often interested in comparing the means or comparing the factor effects. For instance, in a one-factor model, it is of interest to estimate the pairwise differences of the means  $\mu_i - \mu_{i'} = \alpha_i - \alpha_{i'}$ ,  $i \neq i'$ . In general, one may be interested in estimating a linear combination of the means  $\theta = \sum c_i \mu_i$ , where  $\{c_i\}$  are known constants. A linear combination  $\theta = \sum c_i \mu_i$  is called a contrast if  $\sum c_i = 0$ . Thus, any pairwise difference of the means is a contrast.

The least squares estimate of  $\theta = \sum c_i \mu_i$  is  $\hat{\theta} = \sum c_i \bar{Y}_{i..}$ . It is fairly easy to see that its mean, variance, and the estimated variance are

$$\mathbb{E}[\hat{\theta}] = \theta, \quad \text{Var}[\hat{\theta}] = \sigma^2 \sum c_i^2/n_i, \quad \text{and } s^2(\hat{\theta}) = \widehat{\text{Var}[\hat{\theta}]} = \text{MSE} \sum c_i^2/n_i.$$

Since  $\hat{\theta}$  is a function of the estimated mean vector, it is independent of SSE and hence of MSE. Under the assumption of normality  $(\hat{\theta} - \theta)/s(\hat{\theta}) \sim t_{n-k}$ , and this result can be used for constructing confidence intervals or for testing hypotheses.

**Example 11.4.6** (ANCOVA With One Factor and One Covariate). Consider the following model with one factor and one covariate

$$Y_{ij} = \mu + \alpha_i + \gamma Z_{ij} + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, a.$$

If  $\gamma$  were known then we could rewrite the model as  $Y_{ij}^* = \mu + \alpha_i + \varepsilon_{ij}$ , where  $Y_{ij}^* = Y_{ij} - \gamma Z_{ij}$ , and the estimates of  $\mu$  and  $\alpha_i$  would be

$$\tilde{\mu} = \bar{Y}_{..} = \bar{Y}_i - \gamma \bar{Z}_i, \quad \text{and} \quad \tilde{\alpha}_i = \bar{Y}_{i.}^* - \bar{Y}_{..}^* = \bar{Y}_{i.} - \bar{Y}_{..} - \gamma (\bar{Z}_{i.} - \bar{Z}_{..}).$$

When  $\gamma$  is unknown (which is usually the case in practice) and it is estimated by  $\hat{\gamma}$ , then we may plug in the estimate of  $\gamma$  in the above expressions in order to obtain the estimates of  $\mu$  and  $\alpha_i$ . It turns out that this reasoning is valid and it will be discussed in a separate section later.

We now outline a simple strategy for obtaining the least squares estimate of  $\gamma$ . Rewriting the ANCOVA model as

$$Y_{ij} = \mu + \alpha_i^* + \gamma \tilde{Z}_{ij} + \varepsilon_{ij},$$

where  $\tilde{Z}_{ij} = Z_{ij} - \bar{Z}_i$  and  $\alpha_i^* = \alpha_i + \gamma \bar{Z}_i$ , the least squares criterion is

$$\begin{aligned} \sum_i \sum_j (Y_{ij} - \mu - \alpha_i^* - \gamma \tilde{Z}_{ij})^2 &= \sum_i \sum_j Y_{ij}^2 + \sum_i \sum_j (\mu + \alpha_i^*)^2 + \sum_i \sum_j (\gamma \tilde{Z}_{ij})^2 \\ &\quad - 2 \sum_i \sum_j Y_{ij}(\mu + \alpha_i^*) - 2 \sum_i \sum_j Y_{ij}(\gamma \tilde{Z}_{ij}), \end{aligned}$$

since the cross-product term involving  $\mu + \alpha_i^*$  and  $\gamma \tilde{Z}_{ij}$  equals 0. This allows for estimation of  $\gamma$  and  $\mu + \alpha_i^*$  separately. Thus

$$\begin{aligned} \hat{\gamma} &= \sum_i \sum_j Y_{ij} \tilde{Z}_{ij} / \sum_i \sum_j \tilde{Z}_{ij}^2, \\ \hat{\mu} &= \bar{Y}_{..} - \hat{\gamma} \bar{Z}_{..}, \quad \hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..} - \hat{\gamma} (\bar{Z}_{i.} - \bar{Z}_{..}), \quad \text{and} \quad \hat{Y}_{ij} = \bar{Y}_{i.} + \hat{\gamma} \tilde{Z}_{ij}. \end{aligned}$$

Details on ANCOVA models appear in a later section.

## 11.5 Estimation Under Linear Restrictions on $\beta$

In Section 11.4 of this chapter, a simple submodel  $\mathbf{Y} = \beta_0 \mathbf{1} + \boldsymbol{\varepsilon}$  of  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  was considered when discussing the concept of  $R^2$ . However, more general submodels can also be considered, and an analogous result on the decomposition of the residual sum of squares for the submodel can be obtained. The details will be given later, but the result is as follows. If  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  is the estimated mean vector for the full model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  and  $\tilde{\mathbf{Y}} = \mathbf{X}\tilde{\boldsymbol{\beta}}$  is the estimated mean vector under a submodel (reduced model) of the full model, where  $\tilde{\boldsymbol{\beta}}$  is the least squares estimate of  $\boldsymbol{\beta}$  in the submodel, then it turns out that

$$\|\mathbf{Y} - \tilde{\mathbf{Y}}\|^2 = \|\hat{\mathbf{Y}} - \tilde{\mathbf{Y}}\|^2 + \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2.$$

Thus if we write  $SSE_F = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$  and  $SSE_R = \|\mathbf{Y} - \tilde{\mathbf{Y}}\|^2$  as the residual sum of squares under the full and reduced models, then

$$SSE_R = \|\hat{\mathbf{Y}} - \tilde{\mathbf{Y}}\|^2 + SSE_F.$$

Analogous to  $R^2$ , we can obtain the proportional reduction (also known as the *coefficient of partial determination*) in the residual sum of squares when going from the submodel to the full model

$$\frac{SSE_R - SSE_F}{SSE_R}.$$

Let us now see when and how the above results hold. Suppose we wish to estimate the mean vector  $\mathbf{X}\beta$  in the Gauss-Markov model under the restriction  $\mathbf{L}^T\beta = \theta_0$ ,  $\theta_0$  known, where  $L$  is a  $p \times m$  matrix of rank  $m \leq p$ . Such a problem usually comes up in hypothesis testing where the null is  $H_0: \mathbf{L}^T\beta = \theta_0$  against the alternative  $H_1: \mathbf{L}^T\beta \neq \theta_0$ . For such a testing problem, we need to obtain the residual sums of squares for the full model and reduced model (ie, the submodel model with the constraint  $\mathbf{L}^T\beta = \theta_0$ ), and then use them to carry out the test which will be described later.

The restricted least squares estimate of  $\beta$  has a complicated expression and some notational simplifications make the arguments clearer. Since the design matrix  $X$  may not have orthogonal columns, it helps to reexpress  $\mathbf{X}\beta$  as  $\mathbf{X}_0\gamma$  so that the columns of  $\mathbf{X}_0$  are orthonormal. Let  $(\mathbf{X}^T\mathbf{X})^{1/2}$  be a symmetric square root of  $\mathbf{X}^T\mathbf{X}$  and define

$$\begin{aligned} \mathbf{X}_0 &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1/2}, \quad \gamma = (\mathbf{X}^T\mathbf{X})^{1/2}\beta, \quad \mathbf{L}_0 = (\mathbf{X}^T\mathbf{X})^{-1/2}\mathbf{L}, \text{ so that} \\ \mathbf{X}\beta &= \mathbf{X}_0\gamma \text{ and } \theta = \mathbf{L}^T\beta = \mathbf{L}_0^T\gamma. \end{aligned}$$

With these notations, let us note the following:

- (i)  $\mathbf{X}_0^T\mathbf{X}_0 = \mathbf{I}$ ,
- (ii) the least squares estimate of  $\gamma$  in the full model is  $\hat{\gamma} = (\mathbf{X}^T\mathbf{X})^{1/2}\hat{\beta} = \mathbf{X}_0^T\mathbf{Y}$ ,
- (iii)  $E[\hat{\gamma}] = \gamma$ ,  $\text{Cov}[\hat{\gamma}] = \sigma^2\mathbf{I}$ ,
- (iv)  $E[\hat{\theta}] = \theta$  and  $\text{Cov}[\hat{\theta}] = \sigma^2\mathbf{L}_0^T\mathbf{L}_0$ , where  $\hat{\theta} = \mathbf{L}^T\hat{\beta} = \mathbf{L}_0^T\hat{\gamma}$ , and
- (v) the least squares estimate of the mean vector  $\mu = \mathbf{X}\beta$  in the full model is  $\hat{\mu} = \mathbf{X}\hat{\beta} = \mathbf{X}_0\hat{\gamma}$ .

For the restricted least squares case, we now find the estimate of  $\mathbf{X}\beta$ . We minimize  $\|\mathbf{Y} - \mathbf{X}\beta\|^2 = \|\mathbf{Y} - \mathbf{X}_0\gamma\|^2$  with respect to  $\gamma$  subject to the constraint  $\mathbf{L}_0^T\gamma = \theta_0$ . The method of Lagrangian multiplier is useful for a constrained optimization problem, and we minimize

$$\|\mathbf{Y} - \mathbf{X}_0\gamma\|^2 + \lambda^T(\mathbf{L}_0^T\gamma - \theta_0)$$

with respect to  $\gamma$  where  $\lambda$  is the Lagrangian multiplier vector. Differentiating the last expression with respect to  $\gamma$  and  $\lambda$ , and equating the derivatives to  $\mathbf{0}$ , we have

$$\tilde{\gamma} = \mathbf{X}_0^T \mathbf{Y} - (1/2) \mathbf{L}_0 \lambda \quad \text{and} \quad \mathbf{L}_0^T \tilde{\gamma} = \theta_0.$$

Premultiplying the first equation with  $\mathbf{L}_0^T$  we get  $\mathbf{L}_0^T \tilde{\gamma} = \mathbf{L}_0^T \mathbf{X}_0^T \mathbf{Y} - (1/2) \mathbf{L}_0^T \mathbf{L}_0 \lambda$ . Since  $\mathbf{L}_0^T \tilde{\gamma} = \theta_0$ , then  $\lambda = 2(\mathbf{L}_0^T \mathbf{L}_0)^{-1} [\mathbf{L}_0^T \mathbf{X}_0^T \mathbf{Y} - \theta_0]$ . Thus we get a solution to the restricted least squares problem as

$$\begin{aligned}\tilde{\gamma} &= \mathbf{X}_0^T \mathbf{Y} - \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\mathbf{L}_0^T \mathbf{X}_0^T \mathbf{Y} - \theta_0) \\ &= \hat{\gamma} - \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\mathbf{L}_0^T \hat{\gamma} - \theta_0), \text{ and} \\ \tilde{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1/2} \tilde{\gamma}.\end{aligned}$$

The estimated mean vector and the vector of residuals  $\tilde{\epsilon}$  for the restricted least squares are

$$\begin{aligned}\mathbf{X} \tilde{\beta} &= \mathbf{X}_0 \tilde{\gamma} = \mathbf{X}_0 \hat{\gamma} - \mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\mathbf{L}_0^T \hat{\gamma} - \theta_0) \\ &= \mathbf{X} \hat{\beta} - \mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0), \\ \tilde{\epsilon} &= \mathbf{Y} - \mathbf{X} \tilde{\beta} = \hat{\epsilon} + \mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0).\end{aligned}$$

Since  $\mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0)$  is in  $\mathcal{M}(\mathbf{X})$  and  $\hat{\epsilon}$  is orthogonal to  $\mathcal{M}(\mathbf{X})$ , we have

$$\begin{aligned}\|\tilde{\epsilon}\|^2 &= \|\hat{\epsilon}\|^2 + \|\mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0)\|^2 \\ &= \|\hat{\epsilon}\|^2 + (\hat{\theta} - \theta_0)^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0).\end{aligned}$$

Therefore,

$$\begin{aligned}SSE_R &= SSE_F + (\hat{\theta} - \theta_0)^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0), \text{ or} \\ SSE_R - SSE_F &= (\hat{\theta} - \theta_0)^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0).\end{aligned}\tag{2}$$

The discussion above leads to the following important result.

**Theorem 11.5.1.** Let  $\tilde{\beta}$  be the least squares estimate of  $\beta$  in the Gauss-Markov model in Eq. (1) under the restriction  $\mathbf{L}^T \beta = \theta_0$ , where  $\theta_0$  is known, and let  $\mathbf{X} \tilde{\beta}$  be the estimated mean vector under the restriction. Let  $\theta = \mathbf{L}^T \beta$ ,  $\mathbf{L}_0 = (\mathbf{X}^T \mathbf{X})^{-1/2} \mathbf{L}$ , and  $\hat{\theta} = \mathbf{L}^T \hat{\beta}$ , where  $\hat{\beta}$  is the unrestricted least squares estimate of  $\beta$ . Denote the residual vectors  $\hat{\epsilon} = \mathbf{Y} - \mathbf{X} \tilde{\beta}$  and  $\tilde{\epsilon} = \mathbf{Y} - \mathbf{X} \tilde{\beta}$ . Then:

- (a)  $\mathbf{X} \tilde{\beta} = \mathbf{X} \hat{\beta} - \mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0)$ ,
- (b)  $\tilde{\epsilon} = \hat{\epsilon} + \mathbf{X}_0 \mathbf{L}_0 (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0)$ ,
- (c)  $\|\tilde{\epsilon}\|^2 = \|\hat{\epsilon}\|^2 + (\hat{\theta} - \theta_0)^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\theta} - \theta_0)$ .

## 11.6 Gauss-Markov Models: Inference

Throughout this section we assume a Gauss-Markov model with normal errors (ie,  $\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 I)$ ). [Theorem 11.5.1](#) can be used for constructing confidence regions and hypotheses testing. If  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$ , where  $\mathbf{L}$  is a  $p \times m$  matrix of rank  $m \leq p$ , then its least square estimate is  $\hat{\boldsymbol{\theta}} = \mathbf{L}^T \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}}$  is the least square estimate of  $\boldsymbol{\beta}$ . Since  $\hat{\boldsymbol{\theta}}$  is a linear function of the estimated mean vector  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ , it is uncorrelated with the residual vector  $\hat{\boldsymbol{\epsilon}} = \hat{\mathbf{Y}} - \mathbf{X}\hat{\boldsymbol{\beta}}$  ([Theorem 11.3.2](#)). Thus, under the assumption of normality,  $\hat{\boldsymbol{\theta}}$  is independent of  $\hat{\boldsymbol{\epsilon}}$  and hence of  $MSE = \|\hat{\mathbf{Y}} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/(n-p)$ . Note that

- (i)  $\hat{\boldsymbol{\theta}} \sim N_m(\boldsymbol{\theta}, \sigma^2 \mathbf{L}_0^T \mathbf{L}_0)$ , where  $\mathbf{L}_0 = (\mathbf{X}^T \mathbf{X})^{-1/2} \mathbf{L}$ , and  

$$(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) / \sigma^2 \sim \chi_m^2,$$
- (ii)  $SSE/\sigma^2 = \|\hat{\mathbf{Y}} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 / \sigma^2 \sim \chi_{n-p}^2$ ,
- (iii)  $\hat{\boldsymbol{\theta}}$  and  $SSE$  are independent.

The  $F$ -ratio

$$\begin{aligned} F &= \frac{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) / (\sigma^2 m)}{\|\hat{\mathbf{Y}} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 / (\sigma^2 (n-p))} \\ &= \frac{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T (\mathbf{L}_0^T \mathbf{L}_0)^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T / m}{MSE} \\ &= (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T [s^2(\hat{\boldsymbol{\theta}})]^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) / m, \text{ where} \\ s^2(\hat{\boldsymbol{\theta}}) &= MSE [\mathbf{L}_0^T \mathbf{L}_0]^{-1}, \end{aligned}$$

has an  $F$ -distribution with  $df = (m, n-p)$ .

If  $\mathbf{L}$  is a vector (ie,  $m = 1$ ), then  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$  is a real number and  $(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})/s(\hat{\boldsymbol{\theta}}) \sim t_{n-p}$ . Then a confidence interval for  $\boldsymbol{\theta}$  with confidence coefficient  $1 - \alpha$  is  $\hat{\boldsymbol{\theta}} \pm t_{n-p, \alpha/2} s(\hat{\boldsymbol{\theta}})$ . Similarly, if we want to test  $H_0: \boldsymbol{\theta} = \boldsymbol{\theta}_0$  against the alternative  $H_1: \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ , then one may reject the null hypothesis if  $|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/s(\hat{\boldsymbol{\theta}})| > t_{n-p, \alpha/2}$ .

If  $\mathbf{L}$  is a matrix of order  $p \times m$  with rank  $m \leq p$  and  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$ , and it is desired to carry out a test  $H_0: \boldsymbol{\theta} = \boldsymbol{\theta}_0$  against  $H_1: \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ , then there are two equivalent ways to express the  $F$ -statistic for this test. One expression involves  $SSE_F$  and  $SSE_R$ , whereas the other involves  $\hat{\boldsymbol{\theta}} = \mathbf{L}^T \hat{\boldsymbol{\beta}}$  (where  $\hat{\boldsymbol{\beta}}$  is the least squares estimate under the full model) and  $s^2(\hat{\boldsymbol{\theta}})$ , the estimate of  $\text{Cov}[\hat{\boldsymbol{\theta}}]$ . For a particular application, one may use the form that is more convenient to obtain the  $F$ -statistic. If  $MSE_F$  is the mean square error under the full model, then the  $F$ -statistic can be written in the two equivalent forms

$$F = \frac{(SSE_R - SSE_F)/m}{MSE_F} \quad (3a)$$

$$= (\hat{\theta} - \theta_0)^T [s^2(\hat{\theta})]^{-1} (\hat{\theta} - \theta_0) / m. \quad (3b)$$

Since  $\hat{\theta} \sim N_m(\theta, \sigma^2 L^T (X^T X)^{-1} L)$ , using [Theorem 11.5.1](#), we have

$$[SSE_R - SSE_F]/\sigma^2 = (\hat{\theta} - \theta_0)^T [L_0^T L_0]^{-1} (\hat{\theta} - \theta_0) / \sigma^2 \sim \chi_m^2(\delta^2),$$

where  $\delta^2 = (1/2)(\theta - \theta_0)^T [L_0^T L_0]^{-1} (\theta - \theta_0) / \sigma^2$ . Since  $\hat{\beta}$  is independent of the residual vector  $Y - X\hat{\beta}$ ,  $SSE_R - SSE_F$ , which is a function of  $\hat{\beta}$ , is independent of  $SSE_F$ , a function of  $Y - X\hat{\beta}$ . Thus  $SSE_F/\sigma^2 \sim \chi_{n-p}^2$ , and under  $H_0$ ,  $(SSE_R - SSE_F)/\sigma^2 \sim \chi_m^2(\delta^2)$ , where  $\delta^2$  is given above. Therefore

$$\begin{aligned} F &= \frac{(SSE_R - SSE_F)/m}{SSE_F/(n-p)} = \frac{(SSE_R - SSE_F)/m}{MSE_F} \\ &= (\hat{\theta} - \theta_0)^T [s^2(\hat{\theta})]^{-1} (\hat{\theta} - \theta_0) / m \sim F_{m,n-p}(\delta^2). \end{aligned}$$

Thus we reject  $H_0: \theta = \theta_0$  in favor of  $H_1: \theta \neq \theta_0$  if the value of the  $F$ -statistic given above is higher than the critical value obtained from the  $F$ -distribution with  $df = (m, n-p)$ .

A few important facts come out from the above discussions:

- (a) Any linear hypothesis about  $\beta$  induces a reduced model.
- (b) There are two alternate ways to derive the  $F$ -statistic for testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .
- (c)  $SSE_R - SSE_F$  depends on  $Y$  only through  $\hat{\theta} = L^T \hat{\beta}$ . Since  $X\hat{\beta}$  is uncorrelated with  $\hat{\epsilon} = Y - X\hat{\beta}$ ,  $\hat{\beta}$  is independent of  $\hat{\epsilon}$ . Consequently,  $SSE_R - SSE_F$  (a function of  $\hat{\beta}$ ) is independent of  $SSE_F$  (a function of  $\hat{\epsilon}$ ).
- (d)  $df(SSE_F) = n-p$ ,  $df(SSE_R) = n-p+m$ ,  $df(SSE_R) - df(SSE_F) = m$ .
- (e) Under  $H_0: \theta = \theta_0$ , the  $F$ -statistic as given in Eqs. (3a) and (3b) has an  $F$ -distribution with  $df = (m, n-p)$ .

We will summarize the above discussion in the following result.

**Theorem 11.6.1.** Assume the Gauss-Markov setup as in Eq. (1) with  $\epsilon \sim N_n(\mathbf{0}, \sigma^2 I)$ . Let  $L$  be a known matrix of order  $p \times m$  with rank  $m \leq p$ , and denote  $\theta = L^T \beta$ . We wish to test  $H_0: \theta = \theta_0$  against  $H_1: \theta \neq \theta_0$ . Let  $\hat{\beta}$  be the least squares estimate of  $\beta$  and  $\hat{\theta} = L^T \hat{\beta}$ . Let  $SSE_F = \|Y - X\hat{\beta}\|^2$  and  $SSE_R = \|Y - X\tilde{\beta}\|^2$ , where  $\tilde{\beta}$  is the least squares estimate of  $\beta$  under the restriction  $L^T \beta = \theta_0$ . Then

- (a)  $\hat{\theta}$  is independent of  $SSE_F$ ,
- (b)  $SSE_F/\sigma^2 \sim \chi_{n-p}^2$  and  $SSE_R - SSE_F$  is independent of  $SSE_F$ ,
- (c)  $(SSE_R - SSE_F)/\sigma^2 \sim \chi_m^2(\delta^2)$  where the noncentrality parameter is  $\delta^2 = (1/2)(\theta - \theta_0)^T [L_0^T L_0]^{-1} (\theta - \theta_0) / \sigma^2$  with  $L_0 = L(X^T X)^{-1/2}$ , and

(d) for the  $F$ -statistic in Eqs. (3a) and (3b),  $F \sim F_{m,n-p}(\delta^2)$ , where  $\delta^2$  is as given in part (c), and  $F \sim F_{m,n-p}$  under  $H_0$ .

**Example 11.6.1** (Deleting a Variable From the Regression). Let us consider a regression model with  $p - 1 = 4$  independent variables. Suppose we wish to find out if variable  $X_1$  should be dropped from the model. This is equivalent to testing  $H_0: \beta_1 = 0$  against  $H_1: \beta_1 \neq 0$ . In this case, we can write  $\mathbf{L}$  as a row vector of length  $p = 5$  whose second element is 1 and the rest are zeros. In such a case, clearly,  $\theta = \mathbf{L}^T \boldsymbol{\beta} = \beta_1$ . We can carry out either a  $t$ -test or an  $F$ -test for this purpose. The reduced model (ie, the model under  $H_0$ ) has  $p-2$  independent variables instead of  $p-1$ . If we call the original model the full model, then we have

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon \text{ (full)} \\ Y &= \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon \text{ (reduced).} \end{aligned}$$

In this case, we may obtain the residual sums of squares for the full and the reduced models and  $df(SSE_F) = n - 5$ ,  $df(SSE_R) = n - 4$ . The  $F$ -statistic for this test

$$F = \frac{SSE_R - SSE_F}{MSE_F}$$

is exactly equal to  $\left[ \hat{\beta}_1 / s(\hat{\beta}_1) \right]^2$ . Thus we may simply use the  $t$ -statistic  $t = \hat{\beta}_1 / s(\hat{\beta}_1)$  for testing  $H_0: \beta_1 = 0$  against  $H_1: \beta_1 \neq 0$  and avoid calculating  $SSE_R$ .

**Example 11.6.2** (Deleting More Than One Variable From the Regression Model). Suppose the set-up is the same as the last example, but now we wish to know if we can delete the first two independent variables from the model. This is equivalent to testing  $H_0: \beta_1 = \beta_2 = 0$  against  $H_1$ : at least one  $\beta_1, \beta_2$  is nonzero. In this case,  $\mathbf{L}$  can be written as a  $p \times 2$  matrix which has all zeros except for 1's at the second element of the first column and at the third element of the second column. We need to test  $H_0: \theta = \theta_0$  against  $H_1: \theta \neq \theta_0$ , where  $\theta_0 = (0, 0)^T$ . The  $F$ -statistic for this test is  $F = (\hat{\theta} - \theta_0)^T [s^2(\hat{\theta})]^{-1} (\hat{\theta} - \theta_0) / 2$ . However, we can derive this same test statistic in a different way. The full and the reduced (under the null) models are

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon \text{ (full),} \\ Y &= \beta_0 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon \text{ (reduced).} \end{aligned}$$

Here,  $df(SSE_F) = n - 5$ ,  $df(SSE_R) = n - 3$ , and we can write

$$F = \frac{(SSE_R - SSE_F)/2}{MSE_F},$$

and  $F \sim F_{2,n-5}$  under  $H_0$ .

**Example 11.6.3** (One-Factor ANOVA Model). Consider a one-factor model with  $k$  levels. If it is desired to test  $H_0: \alpha_1 = \dots = \alpha_k = 0$  vs  $H_1$ : not all  $\alpha_i$ 's are 0, then the full and the reduced models are

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \text{ (full),} \quad Y_{ij} = \mu + \varepsilon_{ij} \text{ (reduced).}$$

In this case,

$$\begin{aligned} SSE_F &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{i\cdot})^2, \quad df(SSE_F) = n - k, \\ SSE_R &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{..})^2, \quad df(SSE_R) = n - 1, \text{ and} \\ SSR_R - SSE_F &= \sum_i \sum_j (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2 = \sum_i n_i \hat{\alpha}_i^2 := SSTR. \end{aligned}$$

The  $F$ -statistic is

$$F = \frac{[SSR_R - SSE_F]/(k-1)}{MSE_F} = \frac{SSTR/(k-1)}{MSE_F} = \frac{MSTR}{MSE},$$

where  $MSE = MSE_F$ . Denoting  $\delta^2 = (1/2) \sum n_i \alpha_i^2 / \sigma^2$ , we get the result  $F \sim F_{k-1, n-k}(\delta^2)$  and  $F \sim F_{k-1, n-k}$  under  $H_0$ .

**Example 11.6.4** (Balanced Two-Factor ANOVA). In the two-factor balanced ANOVA model, we have seen in [Example 11.4.4](#) that SST0, the total sum of squares, can be decomposed as the sum of SSA, SSB, and SSAB. If we want to test  $H_0: (\alpha\beta)_{ij} = 0$  for all  $i$  and  $j$ , vs  $H_1$ : at least one  $(\alpha\beta)_{ij}$  is not zero, the full and reduced models are

$$\begin{aligned} Y_{ijk} &= \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk} \text{ (full),} \\ Y_{ijk} &= \mu + \alpha_i + \beta_j + \varepsilon_{ijk} \text{ (reduced).} \end{aligned}$$

Estimated mean values for the full and reduced models are

$$\begin{aligned} \hat{Y}_{ijk} &= \bar{Y}_{ij\cdot} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + (\widehat{\alpha\beta})_{ij} \\ \tilde{Y}_{ijk} &= \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j, \end{aligned}$$

where  $\hat{\mu}$ ,  $\hat{\alpha}_i$ , and  $\hat{\beta}_j$  are as given in [Example 11.4.4](#). The sums of squares for the full and the reduced models are

$$\begin{aligned} SSE_F &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij\cdot})^2, \quad df(SSE_F) = n - ab, \\ SSE_R &= \sum_i \sum_j \sum_k (Y_{ijk} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j)^2, \quad df(SSE_R) = n - (a + b - 1), \\ SSE_R - SSE_F &= \sum_i \sum_j \sum_k (\bar{Y}_{ij\cdot} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j)^2 \\ &= n_0 \sum_i \sum_j (\widehat{\alpha\beta})_{ij}^2 = SSAB. \end{aligned}$$

Thus the  $F$ -statistic for this testing problem is

$$F = \frac{SSAB/(ab - a - b + 1)}{MSE_F} = \frac{MSAB}{MSE},$$

where  $MSE$  denotes  $MSE_F$ . Under  $H_0$ ,  $F \sim F_{(a-1)(b-1), n-ab}$ .

Even though it may not be a standard practice to test for the main effects in the presence of interactions, we can still formulate the statistical problem and describe the test statistic. Let  $H_0: \alpha_1 = \dots = \alpha_a = 0$  vs  $H_1: \text{not all } \alpha_i\text{'s are zero}$ . Then the full and reduced models are

$$\begin{aligned} Y_{ijk} &= \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk} \text{ (full),} \\ Y_{ijk} &= \mu + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk} \text{ (reduced).} \end{aligned}$$

Estimated mean values for the full and reduced models, and their sums of squares are

$$\begin{aligned} \hat{Y}_{ijk} &= \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + (\widehat{\alpha\beta})_{ij}, \quad \tilde{Y}_{ijk} = \hat{\mu} + \hat{\beta}_j + (\widehat{\alpha\beta})_{ij}, \\ SSE_F &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij.})^2, \quad df(SSE_F) = n - ab, \\ SSE_R &= \sum_i \sum_j \sum_k (Y_{ijk} - \hat{\mu} - \hat{\beta}_j - (\widehat{\alpha\beta})_{ij})^2, \quad df(SSE_R) = n - [1 + a(b - 1)], \\ SSE_R - SSE_F &= SSA, \end{aligned}$$

and the  $F$ -statistic is

$$F = \frac{SSA/(a - 1)}{MSE_F} = \frac{MSA}{MSE},$$

where  $MSE = MSE_F$ . Under  $H_0$ ,  $F \sim F_{a-1, n-ab}$ .

**Example 11.6.5** (Two-Factor ANOVA: One Observation per Cell). Consider a two-factor ANOVA model as in the previous example with  $n_{ij} = 1$  for all  $i$  and  $j$ . If  $Y_{ij}$  is the response when factor  $A$  is at level  $i$  and factor  $B$  is at level  $j$ , then a model without interactions is

$$Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad j = 1, \dots, b, \quad i = 1, \dots, a,$$

where  $\sum \alpha_i = 0$ ,  $\sum \beta_j = 0$ , and  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ . Estimates of  $\mu$ ,  $\alpha_i$ , and  $\beta_j$  are exactly the same as before (ie,  $\hat{\mu} = \bar{Y}_{..}$ ,  $\hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..}$ , and  $\hat{\beta}_j = \bar{Y}_{.j} - \bar{Y}_{..}$ ). In this case, the fitted  $Y$  values, residuals, and the residual sum of squares are

$$\begin{aligned} \hat{Y}_{ij} &= \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j, \quad \hat{\varepsilon}_{ij} = Y_{ij} - (\hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j), \\ SSE &= \sum_i \sum_j (Y_{ij} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j)^2, \text{ and} \\ df(SSE) &= ab - (a + b - 1) = (a - 1)(b - 1), \text{ and} \\ MSE &= SSE/[(a - 1)(b - 1)]. \end{aligned}$$

One can then carry out inferences on  $\{\alpha_i\}$  and  $\{\beta_j\}$  such as construction of simultaneous confidence intervals or tests such as  $H_0: \alpha_1 = \dots = \alpha_a = 0$  vs  $H_1: \text{not all } \alpha_i \text{ are 0}$ .

Now if it is desired to investigate if the interaction effects are present, one may think of the usual model in the two-factor case

$$Y_{ij} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ij}, \quad j = 1, \dots, b, \quad i = 1, \dots, a,$$

with the usual constraints on  $\{\alpha_i\}$ ,  $\{\beta_j\}$ , and  $\{(\alpha\beta)_{ij}\}$ . For this model, the number of unknown parameters (excluding  $\sigma^2$ ) is  $n = ab$ , the estimated mean is  $\hat{Y}_{ij} = Y_{ij}$  and the

residuals are  $Y_{ij} - \hat{Y}_{ij} = 0$ . Thus there is no way to estimate  $\sigma^2$  as there are too many interaction parameters, and we cannot use this model to determine if the interaction effects are present. One way to approach this issue is to consider a more restrictive type of the interaction of the form  $\theta\alpha_i\beta_j$ ,  $\theta$  real. This leads to the consideration of Tukey's interaction model

$$Y_{ij} = \mu + \alpha_i + \beta_j + \theta\alpha_i\beta_j + \varepsilon_{ij}, \quad j = 1, \dots, b, \quad i = 1, \dots, a,$$

where  $\sum \alpha_i = 0$ ,  $\sum \beta_j = 0$ , and  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ . For this model, absence or presence of interaction effects can be judged by testing  $H_0: \theta = 0$  against  $H_1: \theta \neq 0$ . This is known as "Tukey's one degree of freedom test for nonadditivity", details of which are described below.

Now Tukey's model is no longer a linear model, but if one estimates  $\hat{\mu}$ ,  $\hat{\alpha}_i$ , and  $\hat{\beta}_j$  as in the additive model, then minimizing  $\sum_i \sum_j (Y_{ij} - \tilde{Y}_{ij} - \theta\hat{\alpha}_i\hat{\beta}_j)^2$ , where  $\tilde{Y}_{ij} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j$ , with respect to  $\theta$  leads to an estimate

$$\hat{\theta} = \frac{\sum_i \sum_j \tilde{\varepsilon}_{ij} \hat{\alpha}_i \hat{\beta}_j}{S_\alpha S_\beta}, \text{ where}$$

$$\tilde{\varepsilon}_{ij} = Y_{ij} - \tilde{Y}_{ij}, \quad S_\alpha = \sum \hat{\alpha}_i^2, \quad \text{and} \quad S_\beta = \sum \hat{\beta}_j^2.$$

Since  $\{\tilde{\varepsilon}_{ij}\}$  are independent of  $\{\tilde{Y}_{ij}\}$ ,  $E[\tilde{\varepsilon}_{ij}|\tilde{Y}] = E[\tilde{\varepsilon}_{ij}] = \theta\alpha_i\beta_j$ , where  $\tilde{Y}$  is the vector of  $\tilde{Y}_{ij}$ 's, and

$$\begin{aligned} E[\hat{\theta}|\tilde{Y}] &= \frac{\sum_i \sum_j E[\tilde{\varepsilon}_{ij}|\tilde{Y}] \hat{\alpha}_i \hat{\beta}_j}{S_\alpha S_\beta} = \frac{\sum_i \sum_j (\theta\alpha_i\beta_j) \hat{\alpha}_i \hat{\beta}_j}{S_\alpha S_\beta} \\ &= \theta \frac{\sum_i \sum_j \alpha_i \beta_j \hat{\alpha}_i \hat{\beta}_j}{S_\alpha S_\beta} := \theta A. \end{aligned}$$

For two sequences of constants  $\{e_i, i = 1, \dots, a\}$  and  $\{f_j, j = 1, \dots, b\}$  satisfying the constraints  $\sum e_i = 0$  and  $\sum f_j = 0$ ,

$$\begin{aligned} \sum_i \sum_j \tilde{\varepsilon}_{ij} e_i f_j &= \sum_i \sum_j Y_{ij} e_i f_j, \text{ and} \\ \text{Var} \left[ \sum_i \sum_j \tilde{\varepsilon}_{ij} e_i f_j \middle| \tilde{Y} \right] &= \text{Var} \left[ \sum_i \sum_j \tilde{\varepsilon}_{ij} e_i f_j \right] \\ &= \text{Var} \left[ \sum_i \sum_j Y_{ij} e_i f_j \right] = \sigma^2 \sum_i \sum_j e_i^2 f_j^2. \end{aligned}$$

Thus

$$\text{Var}[\hat{\theta}|\tilde{\mathbf{Y}}] = \sigma^2/[S_\alpha S_\beta].$$

In order to carry out a test on  $\theta$ , we need an estimate of  $\sigma^2$ . If we write  $\hat{Y}_{ij} = \tilde{Y}_{ij} + \hat{\theta}\hat{\alpha}_i\hat{\beta}_j$  and  $\hat{\varepsilon}_{ij} = Y_{ij} - \hat{Y}_{ij} = \tilde{\varepsilon}_{ij} - \hat{\theta}\hat{\alpha}_i\hat{\beta}_j$ , then a simple calculation will show that

$$SSE_R = \sum_i \sum_j \tilde{\varepsilon}_{ij}^2 = \sum_i \sum_j \hat{\varepsilon}_{ij}^2 + \hat{\theta}^2 S_\alpha S_\beta = SSE_F + \hat{\theta}^2 S_\alpha S_\beta.$$

Conditionally on  $\tilde{\mathbf{Y}}$ ,

- (i)  $[SSE_R - SSE_F]/\sigma^2 = \hat{\theta}^2 S_\alpha S_\beta/\sigma^2 \sim \chi_1^2(\delta_1^2)$ ,  
where  $\delta_1^2 = (1/2)\{\mathbb{E}[\tilde{\theta}|\tilde{\mathbf{Y}}]\}^2/\sigma^2 = (1/2)\theta^2 A^2/\sigma^2$ ,
- (ii)  $SSE_R/\sigma^2 \sim \chi_{(a-1)(b-1)}^2(\delta^2)$ , where  $\delta^2 = \theta^2 \sum_j \sum_i \alpha_i^2 \beta_j^2/\sigma^2$ , and
- (iii)  $SSE_F$  is nonnegative.

Hence by an application of [Lemma B.7.2](#) in [Appendix B](#), we can conclude that, conditionally on  $\tilde{\mathbf{Y}}$ ,  $SSE_F/\sigma^2 \sim \chi_{(a-1)(b-1)-1}^2(\delta^2 - \delta_1^2)$  and  $SSE_F/\sigma^2$  is independent of  $[SSE_R - SSE_F]/\sigma^2$ . This allows us to construct a test for  $H_0: \theta = 0$  vs  $H_1: \theta \neq 0$  using the  $F$ -statistic

$$F = \frac{SSE_R - SSE_F}{MSE_F} = \frac{\hat{\theta}^2 S_\alpha S_\beta}{MSE_F}, \text{ where}$$

$$MSE_F = \left[ \sum_i \sum_j \tilde{\varepsilon}_{ij}^2 - \hat{\theta}^2 S_\alpha S_\beta \right] / [(a-1)(b-1)-1],$$

and  $F \sim F_{1,(a-1)(b-1)-1}$  under  $H_0$ .

**Example 11.6.6** (Nested ANOVA). Let us begin with a simple example. The school superintendent has asked every school in a town to try a pilot training program in order to improve the quantitative aptitude of the students. The town has  $a$  schools and each school has its own  $b$  designated teachers for this training program. A random sample of  $n_{ij}$  students is chosen in the  $i$ th school (factor  $A$ ),  $i = 1, \dots, a$ , and assigned to the  $j$ th teacher in that school,  $j = 1, \dots, b$ , and after 6 months of training, the students are given a standardized test to evaluate their performances. Note that the teachers in different schools are entirely different and thus the teacher effect (factor  $B$ ) is nested in the school (factor  $A$ ). If  $Y_{ijk}$  is the score of the  $k$ th student assigned to teacher  $j$  in the  $i$ th school, then a reasonable model is

$$Y_{ijk} = \mu + \alpha_i + \beta_{j(i)} + \varepsilon_{ijk}, \quad k = 1, \dots, n_{ij}, \quad j = 1, \dots, b, \quad i = 1, \dots, a,$$

where  $\{\alpha_i\}$  are the factor  $A$  (school) effects,  $\{\beta_{j(i)}\}$  are the factor  $B$  (teacher) effects nested in factor  $A$ , and  $\{\varepsilon_{ijk}\}$  are iid  $N(0, \sigma^2)$  errors. This is an example of a simple *nested ANOVA model*. For this model, it is assumed that

$$(a) \sum \alpha_i = 0, \quad (b) \sum_{j=1}^b \beta_{j(i)} = 0 \quad \text{for each } i = 1, \dots, a.$$

Let  $\bar{Y}_{...}$ ,  $\bar{Y}_{i..}$ , etc., be as in the two-factor ANOVA model. Then estimates of  $\mu$ ,  $\alpha_i$ ,  $\beta_{j(i)}$ , fitted  $Y$  values, and the residuals are

$$\begin{aligned}\hat{\mu} &= (ab)^{-1} \sum_i \sum_j \bar{Y}_{ij.}, \quad \hat{\alpha}_i = b^{-1} \sum_j \bar{Y}_{ij.} - \hat{\mu}, \quad \hat{\beta}_{j(i)} = \bar{Y}_{ij.} - \hat{\mu} - \hat{\alpha}_i, \\ \hat{Y}_{ijk} &= \bar{Y}_{ij.}, \quad \hat{\varepsilon}_{ijk} = Y_{ijk} - \bar{Y}_{ij.}, \quad \text{and} \\ SSE &= \sum_i \sum_j \sum_k \hat{\varepsilon}_{ijk}^2.\end{aligned}$$

As usual,  $SSE/\sigma^2 \sim \chi^2_{n-ab}$ , and  $MSE = SSE/(n-ab)$  is an unbiased estimate of  $\sigma^2$ , where  $n$  is the total number of observations. One can obtain variances of  $\hat{\mu}$ ,  $\{\hat{\alpha}_i\}$ , and  $\{\hat{\beta}_{j(i)}\}$ , but they are a bit cumbersome in the unbalanced case. In the balanced case (ie,  $n_{ij} = n_0$  for all  $i$  and  $j$ ), the expressions of  $\hat{\mu}$ ,  $\{\hat{\alpha}_i\}$ , and  $\{\hat{\beta}_{j(i)}\}$  are simpler

$$\hat{\mu} = \bar{Y}_{...}, \quad \hat{\alpha}_i = \bar{Y}_{i..} - \bar{Y}_{...}, \quad \hat{\beta}_{j(i)} = \bar{Y}_{ij.} - \hat{\mu} - \hat{\alpha}_i = \bar{Y}_{ij.} - \bar{Y}_{i..},$$

and

$$\begin{aligned}\text{Var}[\hat{\mu}] &= \text{Var}[\bar{Y}_{...}] = \sigma^2/n, \\ \text{Var}[\hat{\alpha}_i] &= \text{Var}[\bar{Y}_{i..} - \bar{Y}_{...}] = \sigma^2(a-1)/n, \\ \text{Var}[\hat{\beta}_{j(i)}] &= \text{Var}[\bar{Y}_{ij.} - \bar{Y}_{i..}] = \sigma^2(b-1)/(n_0 b).\end{aligned}$$

Since  $\hat{\mu}$ ,  $\{\hat{\alpha}_i\}$ , and  $\{\hat{\beta}_{j(i)}\}$  are functions of the fitted mean vector  $\hat{Y}$ , they are independent of  $SSE$  and hence of  $MSE$ , and this fact can be used to carry out inferences such as hypotheses testing and construction of confidence intervals.

If we wish to test the hypothesis of no teacher effect (ie, test  $H_0: \beta_{j(i)} = 0$  for all  $j$  and  $i$ ) against the alternative  $H_1$ : not all  $\beta_{j(i)}$  are zero, then the full and reduced models and fitted  $Y$  values are

$$\begin{aligned}Y_{ijk} &= \mu + \alpha_i + \beta_{j(i)} + \varepsilon_{ijk}, \quad \hat{Y}_{ijk} = \bar{Y}_{ij.} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_{j(i)} \text{ (full),} \\ Y_{ijk} &= \mu + \alpha_i + \varepsilon_{ijk}, \quad \tilde{Y}_{ijk} = \hat{\mu} + \hat{\alpha}_i \text{ (reduced).}\end{aligned}$$

The residual sums of squares are

$$\begin{aligned}SSE_F &= \sum_i \sum_j \sum_k \left( Y_{ijk} - \bar{Y}_{ij.} \right)^2, \quad df(SSE_F) = n - ab, \\ SSE_R &= \sum_i \sum_j \sum_k \left( Y_{ijk} - \hat{\mu} - \hat{\alpha}_i \right)^2, \quad df(SSE_R) = n - a, \\ SSE_R - SSE_F &= \sum_i \sum_j n_{ij} \hat{\beta}_{j(i)}^2 := SSB(A).\end{aligned}$$

The quantity  $SSB(A)$  is the sum of squares due to teachers (factor  $B$ ) nested in school (factor  $A$ ). The test statistic is

$$F = \frac{SSB(A)/(ab - a)}{MSE_F} = \frac{MSB(A)}{MSE},$$

where  $MSB(A) = SSB(A)/(ab - a)$  and  $MSE = MSE_F$ . Under  $H_0$ ,  $F \sim F_{ab-a, n-ab}$ .

In order to test the hypothesis of no school effect (ie,  $H_0: \alpha_1 = \dots = \alpha_a = 0$  vs  $H_1$ : at least one  $\alpha_i$  is not zero), the full and reduced models along with the fitted values are

$$\begin{aligned} Y_{ijk} &= \mu + \alpha_i + \beta_{j(i)} + \varepsilon_{ijk}, \quad \hat{Y}_{ijk} = \bar{Y}_{ij\cdot} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_{j(i)} \text{ (full),} \\ Y_{ijk} &= \mu + \beta_{j(i)} + \varepsilon_{ijk}, \quad \hat{Y}_{ijk} = \hat{\mu} + \hat{\beta}_{j(i)} \text{ (reduced).} \end{aligned}$$

The residual sum of squares are

$$\begin{aligned} SSE_F &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij\cdot})^2, \quad df(SSE_F) = n - ab, \\ SSE_R &= \sum_i \sum_j \sum_k (Y_{ijk} - \hat{\mu} - \hat{\beta}_{j(i)})^2, \\ df(SSE_R) &= n - 1 - a(b - 1), \text{ and} \\ SSE_R - SSE_F &= \sum_i \sum_j n_{ij} \hat{\alpha}_i^2 := SSA. \end{aligned}$$

Thus the  $F$ -statistic

$$F = \frac{SSA/(a-1)}{MSE_F} = \frac{MSA}{MSE},$$

where  $MSE = MSE_F$ , follows an  $F$ -distribution with  $df(a-1, n-ab)$  under  $H_0$ .

For the balanced case, the total sum of squares admits the following decomposition

$$SSTO = SSA + SSB(A) + SSE.$$

### 11.6.1 Simultaneous Inference

We now address the issue of constructing simultaneous confidence intervals for  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$  and its linear functions, where  $\mathbf{L}$  is a matrix of order  $p \times m$  of rank  $m \leq p$ . The main ingredients are the basic distributional results on the least squares estimate  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  and  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$ , as outlined in this section. Recall that the estimate of the error variance  $\sigma^2$  in the Gauss-Markov model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  is given by  $MSE = \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/(n-p)$ . Thus for a real-valued parameter  $\theta = \mathbf{L}^T \boldsymbol{\beta}$ ,  $(\hat{\theta} - \theta)/s(\hat{\theta}) \sim t_{n-p}$  where  $\hat{\theta} = \mathbf{L}^T \hat{\boldsymbol{\beta}}$  and  $s^2(\hat{\theta}) = MSE \left\{ \mathbf{L}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L} \right\}^{-1}$ . Hence  $\hat{\theta} \pm t_{n-p, \alpha/2} s(\hat{\theta})$  is a confidence interval for  $\theta$  with confidence coefficient  $1 - \alpha$ . Now if  $\boldsymbol{\theta}$  is vector valued (ie,  $\mathbf{L}$  is a  $p \times m$  matrix), then there are different methods for constructing simultaneous confidence intervals for  $\theta_1, \dots, \theta_m$ , the components of  $\boldsymbol{\theta}$ , and for the linear combinations of  $\boldsymbol{\theta}$ . In the literature, there are many methods for simultaneous inference, but the discussion below will be for only three

well-known methods: Bonferroni, Scheffé, and Tukey. Even though we only discuss the problem of constructing simultaneous confidence intervals, these can also be employed for simultaneous hypotheses testing.

### Bonferroni Method

According to the Bonferroni method, simultaneous confidence intervals with a family confidence coefficient of at least  $1 - \alpha$  are given by

$$\theta_j: \hat{\theta}_j \pm Bs(\hat{\theta}_j), \quad j = 1, \dots, m, \quad \text{with } B = t_{n-p,\alpha/(2m)},$$

where  $s^2(\hat{\theta}_j)$  is the  $j$ th diagonal element of the matrix  $s^2(\hat{\theta})$ . Even though this method is valid for any  $m$ , its usefulness is questionable when  $m$  is not small since the multiplier  $t_{n-p,\alpha/(2m)}$  associated with the confidence intervals increases as  $m$  increases. Mathematically, the multiplier converges to  $\infty$  as  $m \rightarrow \infty$ . In reality, the Bonferroni method is an inefficient method for constructing simultaneous confidence intervals when  $m$  is larger than 3 or 4. The Scheffé procedure is more appropriate when  $m$  is large. Before describing the Scheffé method, let us briefly see why the Bonferroni method leads to a simultaneous confidence of at least  $1 - \alpha$ . Let  $A_j$  denote the random event  $\{\theta_j: \theta_j \in [\hat{\theta}_j - Bs(\hat{\theta}_j), \hat{\theta}_j + Bs(\hat{\theta}_j)]\}$ ,  $j = 1, \dots, m$ . Now  $P(A_j^c) = \alpha/m$  and

$$P\left(\bigcap_{j=1}^m A_j\right) = 1 - P\left(\bigcup_{j=1}^m A_j^c\right) \geq 1 - \sum_{j=1}^m P(A_j^c) = 1 - \alpha.$$

It shows that the probability that  $\theta_j$  is inside  $\hat{\theta}_j \pm Bs(\hat{\theta}_j)$  for all  $j = 1, \dots, m$ , is at least  $1 - \alpha$ , which justifies the Bonferroni approach. This leads to the following lemma.

**Lemma 11.6.1.** *The confidence intervals  $\hat{\theta}_j \pm Bs(\hat{\theta}_j)$  for  $\theta_j$ ,  $j = 1, \dots, m$ , have a simultaneous confidence of at least  $1 - \alpha$ , where  $B = t_{n-p,\alpha/(2m)}$ .*

### Scheffé Method

This method obtains simultaneous confidence intervals for all linear combinations  $\mathbf{a}^T \boldsymbol{\theta}$ ,  $\mathbf{a} \in \mathbb{R}^m$ . Note that the least squares estimate of  $\mathbf{a}^T \boldsymbol{\theta}$  is  $\mathbf{a}^T \hat{\boldsymbol{\theta}} = \mathbf{a}^T L^T \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}}$  is the least squares estimate of  $\boldsymbol{\beta}$ . It is also clear that  $E[\mathbf{a}^T \hat{\boldsymbol{\theta}}] = \mathbf{a}^T \boldsymbol{\theta}$ ,  $\text{Var}[\mathbf{a}^T \hat{\boldsymbol{\theta}}] = \mathbf{a}^T \text{Cov}[\hat{\boldsymbol{\theta}}] \mathbf{a}$ , and an estimate of  $\text{Var}[\mathbf{a}^T \hat{\boldsymbol{\theta}}]$  is given by  $s^2(\mathbf{a}^T \hat{\boldsymbol{\theta}}) = \mathbf{a}^T s^2(\hat{\boldsymbol{\theta}}) \mathbf{a}$ .

This method states that simultaneous confidence intervals for all linear combinations of  $\boldsymbol{\theta}$  with an overall confidence level of  $1 - \alpha$  are given by

$$\mathbf{a}^T \boldsymbol{\theta}: \mathbf{a}^T \hat{\boldsymbol{\theta}} \pm Ss(\mathbf{a}^T \hat{\boldsymbol{\theta}}), \quad \mathbf{a} \in \mathbb{R}^m, \quad \text{where } S = \sqrt{m F_{m,n-p,\alpha}}.$$

The Scheffé method is closely related to a procedure for obtaining a confidence region for  $\theta$ , known as the confidence ellipsoid, which is given by

$$A = \left\{ \theta \in \mathbb{R}^m : (\theta - \hat{\theta})^T [s^2(\hat{\theta})]^{-1} (\theta - \hat{\theta}) / m \leq F_{m,n-p,\alpha} \right\}.$$

Note that  $P_\theta[A] = 1 - \alpha$  and thus the confidence ellipsoid provides a confidence region for  $\theta$  with confidence  $1 - \alpha$ . However, this is not a useful method in practice as it is not possible to visualize the region when  $m > 3$ .

There is a connection between Scheffé's simultaneous confidence intervals and the confidence ellipsoid, and it is through the following equality

$$\sup_{\mathbf{a} \in \mathbb{R}^m} (\mathbf{a}^T \hat{\theta} - \mathbf{a}^T \theta)^2 / s^2(\mathbf{a}^T \hat{\theta}) = (\theta - \hat{\theta})^T [s^2(\hat{\theta})]^{-1} (\theta - \hat{\theta}),$$

which basically follows from the Cauchy-Schwarz inequality (Section B.1) since

$$(\mathbf{a}^T \hat{\theta} - \mathbf{a}^T \theta)^2 / s^2(\mathbf{a}^T \hat{\theta}) = [\mathbf{a}^T (\hat{\theta} - \theta)]^2 / [\mathbf{a}^T s^2(\hat{\theta}) \mathbf{a}].$$

Thus we have

$$\begin{aligned} 1 - \alpha &= P_\theta \left[ (\hat{\theta} - \theta)^T [s^2(\hat{\theta})]^{-1} (\hat{\theta} - \theta) / m \leq F_{m,n-p,\alpha} \right] \\ &= P_\theta \left[ \sup_{\mathbf{a} \in \mathbb{R}^m} (\mathbf{a}^T \hat{\theta} - \mathbf{a}^T \theta)^2 / [\mathbf{a}^T s^2(\hat{\theta}) \mathbf{a}] \leq S^2 \right] \\ &= P_\theta \left[ (\mathbf{a}^T \hat{\theta} - \mathbf{a}^T \theta)^2 / [\mathbf{a}^T s^2(\hat{\theta}) \mathbf{a}] \leq S^2 \text{ for all } \mathbf{a} \in \mathbb{R}^m \right] \\ &= P_\theta \left[ \mathbf{a}^T \hat{\theta} - Ss(\mathbf{a}^T \hat{\theta}) \leq \mathbf{a}^T \theta \leq \mathbf{a}^T \hat{\theta} + Ss(\mathbf{a}^T \hat{\theta}) \text{ for all } \mathbf{a} \in \mathbb{R}^m \right]. \end{aligned}$$

This justifies the validity of the Scheffé method and we have the following result.

**Theorem 11.6.2.** *Let  $\theta$  be a linear function of  $\beta$  of the form  $\mathbf{L}^T \beta$  where  $\mathbf{L}$  is a  $p \times m$  matrix of rank  $m \leq p$ . Then simultaneous confidence intervals for all linear combinations of  $\theta$  with an overall confidence coefficient of  $1 - \alpha$  are*

$$\mathbf{a}^T \theta : \mathbf{a}^T \hat{\theta} \pm Ss(\mathbf{a}^T \hat{\theta}), \quad \mathbf{a} \in \mathbb{R}^m, \text{ where } S = \sqrt{mF_{m,n-p,\alpha}}.$$

## Tukey Method

The basic argument behind Tukey's method is given below and it will be clear from an example (given below) how the method can be used. Let  $W_i \sim N(\theta_i, \tau^2)$ ,  $i = 1, \dots, t$ , be independent random variables and let  $S^2$  be an unbiased estimate of  $\tau^2$  such that

- (i)  $S^2$  is independent of  $W_1, \dots, W_t$ ,
- (ii)  $vS^2/\tau^2 \sim \chi_v^2$ .

Consider the following rv, called the studentized range variable,

$$Q = \left\{ \max_i(W_i - \theta_i) - \min_i(W_i - \theta_i) \right\} / S.$$

The distribution of  $Q$  is known as the studentized range distribution (denoted by  $Q(t, v)$ ) and it is available in statistical packages. Many statistics textbooks also have the table of this distribution. Tukey simultaneous confidence intervals with family confidence coefficient  $1 - \alpha$  for all pairwise differences  $\theta_i - \theta_j$  is

$$W_i - W_j \pm Q_{t, v, \alpha} S,$$

where  $Q_{t, v, \alpha}$  is the  $(1 - \alpha)$ -quantile of the studentized range distribution with degrees of freedom  $(t, v)$ .

### Application of Tukey's Method to One-Factor Balanced ANOVA Models

Let  $Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$ ,  $j = 1, \dots, n_0$ ,  $i = 1, \dots, k$ , where  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$  and let  $\mu_i = \mu + \alpha_i$ ,  $i = 1, \dots, k$ . Suppose that we want to construct simultaneous confidence intervals for  $\mu_i - \mu_{i'} = \alpha_i - \alpha_{i'}$ ,  $i \neq i'$ , with a family confidence coefficient of  $1 - \alpha$ . Take  $W_i = \sqrt{n_0} \bar{Y}_i$ ,  $\theta_i = \sqrt{n_0} \mu_i$ ,  $S^2 = MSE$ ,  $\tau^2 = \sigma^2$ ,  $t = k$ , and  $v = n - k$ , where  $n = kn_0$  is the total number of observations. Then simultaneous confidence intervals for  $\theta_i - \theta_{i'} = \sqrt{n_0}(\alpha_i - \alpha_{i'})$ ,  $i \neq i'$ , with family confidence coefficient  $1 - \alpha$  are given by

$$\sqrt{n_0}(\bar{Y}_{i \cdot} - \bar{Y}_{i' \cdot}) \pm Q_{k, n-k, \alpha} MSE.$$

Consequently, simultaneous confidence intervals for  $\theta_i - \theta_{i'}$ ,  $i \neq i'$ , are

$$(\bar{Y}_{i \cdot} - \bar{Y}_{i' \cdot}) \pm T S(\bar{Y}_{i \cdot} - \bar{Y}_{i' \cdot}),$$

where  $T = Q_{k, n-k, \alpha} / \sqrt{2}$ .

#### 11.6.2 Prediction Intervals

If we rewrite the Gauss-Markov model as

$$Y_i = \beta^T \mathbf{x}_i + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\mathbf{x}_i^T$  is the  $i$ th row of the design matrix, it is sometimes of interest to predict  $Y_0$  at a vector  $\mathbf{x}_0$  of values in  $\mathbb{R}^p$  where  $Y_0 = \beta^T \mathbf{x}_0 + \varepsilon_0$  and  $\varepsilon_0 \sim N(0, \sigma^2)$  is independent of the observation vector  $\mathbf{Y}$ . If  $\beta$  were known, the best predictor of  $Y_0$  would be  $\beta^T \mathbf{x}_0$ . Usually,  $\beta$  is unknown and the predicted value of  $Y_0$  is  $\hat{Y}_0 = \hat{\beta}^T \mathbf{x}_0$ , where  $\hat{\beta}$  is the least squares estimate of  $\beta$ . Noting that  $\varepsilon_0$  is independent of  $\hat{\beta}$  and hence of  $\hat{Y}_0$ , we have

$$\begin{aligned} E[\hat{Y}_0 - Y_0] &= E[\hat{\beta}^T \mathbf{x}_0 - \beta^T \mathbf{x}_0 - \varepsilon_0] = 0, \text{ and} \\ \text{Var}[\hat{Y}_0 - Y_0] &= \text{Var}[\hat{\beta}^T \mathbf{x}_0 - \beta^T \mathbf{x}_0 - \varepsilon_0] = \text{Var}[\hat{\beta}^T \mathbf{x}_0 - \beta^T \mathbf{x}_0] + \text{Var}[\varepsilon_0] \\ &= \mathbf{x}_0^T \text{Cov}[\hat{\beta}] \mathbf{x}_0 + \sigma^2 = \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 + \sigma^2. \end{aligned}$$

Thus an estimate of  $\text{Var}[\hat{Y}_0 - Y_0]$  is

$$s^2(\hat{Y}_0 - Y_0) = \text{MSE } \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 + \text{MSE} := s^2(\text{pred}).$$

Since  $\hat{Y}_0 - Y_0$  is normally distributed with mean 0 and variance given above, and is independent of  $s^2(\hat{Y}_0 - Y_0)$ , the random variable  $(\hat{Y}_0 - Y_0)/s(\text{pred}) \sim t_{n-p}$ . Thus a prediction interval for  $Y_0$  with confidence coefficient  $1 - \alpha$  is given by  $\hat{Y}_0 \pm t_{n-p,\alpha/2}s(\text{pred})$ .

Simultaneous prediction intervals for  $m$  different  $Y$  values at  $k$  different  $\mathbf{x}$  vectors can be constructed using the Bonferroni method by taking the multiplier associated with the prediction intervals to be equal to  $B = t_{n-p,\alpha/(2m)}$ .

## 11.7 Analysis of Covariance

As we have discussed in Remark 11.2.2, the ANCOVA model can be written as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}, \quad (4)$$

where  $\mathbf{Y}$  is  $n \times 1$  vector of observations,  $\mathbf{X}$  is a known  $n \times p$  matrix of rank  $p$ ,  $\mathbf{Z}$  is a known  $n \times q$  matrix of rank  $q$ , and  $\boldsymbol{\varepsilon}$  is  $N_n(\mathbf{0}, \sigma^2 I)$ . We further assume that  $\mathcal{M}(\mathbf{X}) \cap \mathcal{M}(\mathbf{Z}) = \{\mathbf{0}\}$ , and hence the rank of the augmented matrix  $[\mathbf{X} \ \mathbf{Z}]$  is  $p + q$ . Here  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  are vectors of unknown parameters to be estimated.

Here we will be concerned with estimation of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , and inference on them. For instance, in Example 11.2.7 we may be interested in testing that diet has no effect on the growth rate (ie,  $H_0: \alpha_1 = \dots = \alpha_k = 0$  vs  $H_1$ : not all  $\alpha_i$ 's equal 0). We may be interested in testing that age has no effect on growth (ie,  $H_0: \gamma_2 = 0$  vs  $H_1: \gamma_2 \neq 0$ ). Or we may be interested in finding out neither initial weight nor age has any effect on the growth rate,  $H_0: \gamma_1 = \gamma_2 = 0$  vs  $H_1$ : not both of  $\gamma_1$  and  $\gamma_2$  are zero.

Here we discuss the following issues:

- (i) estimation of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , and their linear functions  $\boldsymbol{\theta} = \mathbf{L}^T \boldsymbol{\beta}$  and  $\boldsymbol{\eta} = \mathbf{M}^T \boldsymbol{\gamma}$ ,
- (ii) test for  $H_0: \boldsymbol{\theta} = \boldsymbol{\theta}_0$  vs  $H_1: \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ ,  $\boldsymbol{\theta}_0$  known, and
- (iii) test for  $H_0: \boldsymbol{\eta} = \boldsymbol{\eta}_0$  vs  $H_1: \boldsymbol{\eta} \neq \boldsymbol{\eta}_0$ ,  $\boldsymbol{\eta}_0$  known.

### 11.7.1 Estimation of $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$

If  $\mathbf{W} = [\mathbf{X} \ \mathbf{Z}]$  and  $\boldsymbol{\theta}$  is the  $(p+q)$ -dim vector obtained by stacking  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  vertically, then the normal equations would be  $\mathbf{W}^T \mathbf{W} \boldsymbol{\theta} = \mathbf{W}^T \mathbf{Y}$ . The matrix  $\mathbf{W}^T \mathbf{W} = \begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{X} & \mathbf{Z}^T \mathbf{Z} \end{pmatrix}$  is not necessarily block diagonal. It would help in obtaining a simple estimate of  $\boldsymbol{\gamma}$  if we could rewrite the model in order to get a block diagonal matrix on the left-hand side of the normal equations. In order to achieve this, we can argue as follows. Let  $\mathbf{Q}_X$

be the orthogonal projection on the column space of  $\mathbf{X}$  (ie,  $\mathbf{Q}_X = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  and  $\tilde{\mathbf{Z}} = (\mathbf{I} - \mathbf{Q}_X)\mathbf{Z}$ ). Then we can rewrite the ANCOVA model as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\delta} + \tilde{\mathbf{Z}}\boldsymbol{\gamma} + \boldsymbol{\varepsilon},$$

where  $\boldsymbol{\delta} = \boldsymbol{\beta} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Z}\boldsymbol{\gamma}$ . Since  $\mathbf{X}^T\tilde{\mathbf{Z}} = \mathbf{0}$ , the least squares method produces the normal equations

$$\begin{pmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\delta} \\ \boldsymbol{\gamma} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{Y} \\ \tilde{\mathbf{Z}}^T\mathbf{Y} \end{pmatrix}.$$

From the second set of equations we get

$$\hat{\boldsymbol{\gamma}} = (\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}})^{-1}\tilde{\mathbf{Z}}^T\mathbf{Y}.$$

From the first set of equations we get

$$\begin{aligned} \mathbf{X}^T\mathbf{Y} &= \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\delta}} = \mathbf{X}^T\mathbf{X}\left(\hat{\boldsymbol{\beta}} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Z}\hat{\boldsymbol{\gamma}}\right) = \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}^T\mathbf{Z}\hat{\boldsymbol{\gamma}}, \text{ and} \\ \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{X}^T(\mathbf{Y} - \mathbf{Z}\hat{\boldsymbol{\gamma}}). \end{aligned}$$

The discussion above leads to the following simple result.

**Lemma 11.7.1.** *The least squares estimates of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  for the ANCOVA model are*

$$(a) \hat{\boldsymbol{\gamma}} = (\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}})^{-1}\tilde{\mathbf{Z}}^T\mathbf{Y}, (b) \hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{Y} - \mathbf{Z}\hat{\boldsymbol{\gamma}}).$$

The following lemma states a few basic results on the least squares estimates and its proof is left to the reader.

**Lemma 11.7.2.** *If  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$  are least squares estimates of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , respectively, then we have*

- (a)  $E[\hat{\boldsymbol{\gamma}}] = \boldsymbol{\gamma}$ ,  $E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ ,
- (b)  $Cov[\hat{\boldsymbol{\gamma}}] = \sigma^2(\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}})^{-1}$ ,
- (c)  $Cov[\hat{\boldsymbol{\beta}}] = \sigma^2\left[(\mathbf{X}^T\mathbf{X})^{-1} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{D}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\right]$ , where  $\mathbf{D} = \mathbf{Z}(\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}})^{-1}\mathbf{Z}^T$ ,
- (d)  $Cov[X\hat{\boldsymbol{\beta}}] = \sigma^2(\mathbf{Q}_X + \mathbf{Q}_X\mathbf{D}\mathbf{Q}_X)$ ,
- (e)  $Cov[\hat{\boldsymbol{\gamma}}, X\hat{\boldsymbol{\beta}}] = -\sigma^2(\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}})^{-1}\mathbf{Z}^T\mathbf{Q}_X$ .

*Remark 11.7.1.* Here is an intuitive way to view the matrix  $\tilde{\mathbf{Z}}$ . Suppose that the columns of  $\mathbf{Z}$  are  $\mathbf{Z}_1, \dots, \mathbf{Z}_q$ . Then the  $i$ th column vector for the matrix  $\tilde{\mathbf{Z}}$  is  $\tilde{\mathbf{Z}}_i = (\mathbf{I} - \mathbf{Q}_X)\mathbf{Z}_i$ . Now for any column of  $\mathbf{Z}$ , say  $\mathbf{Z}_1$ , we can view  $\mathbf{Q}_X\mathbf{Z}_1$  as the vector of fitted values when fitting the model  $\mathbf{Z}_1 = \mathbf{X}\boldsymbol{\beta} + \text{error}$ , and hence the vector of residuals is  $(\mathbf{I} - \mathbf{Q}_X)\mathbf{Z}_1$ . Similar interpretation holds for all the columns of  $\tilde{\mathbf{Z}}$ .

### 11.7.2 Residual Sum of Squares

The vectors of fitted values and the residuals for the ANCOVA model are

$$\begin{aligned}\hat{Y} &= X\hat{\beta} + Z\hat{\gamma} = Q_X(Y - Z\hat{\gamma}) + Z\hat{\gamma} = Q_X Y + \tilde{Z}\hat{\gamma}, \\ \hat{\epsilon} &= Y - \hat{Y} = (I - Q_X)Y - \tilde{Z}\hat{\gamma}.\end{aligned}$$

Since  $\hat{\epsilon}$  is orthogonal to any linear combination of  $X$  and  $Z$ , and hence to  $\tilde{Z}\hat{\gamma}$ , we have

$$\begin{aligned}\|(I - Q_X)Y\|^2 &= \|\hat{\epsilon} + \tilde{Z}\hat{\gamma}\|^2 = \|\hat{\epsilon}\|^2 + \|\tilde{Z}\hat{\gamma}\|^2, \text{ and hence} \\ SSE &= \|\hat{\epsilon}\|^2 = \|(I - Q_X)Y\|^2 - \|\tilde{Z}\hat{\gamma}\|^2.\end{aligned}$$

Since the rank of the matrix  $[X \ Z]$  is  $p + q$  and  $df(SSE) = n - p - q$ , an unbiased estimate of  $\sigma^2$  is

$$MSE = SSE/(n - p - q).$$

### 11.7.3 Inference for $\gamma$ and $\beta$

We first discuss inference on  $\gamma$  and its linear functions.

#### *Inference for $\gamma$*

If we want to estimate  $\eta = M^T\gamma$ , where  $M$  is a  $q \times s$  matrix of rank  $s \leq q$ , then the least squares estimate of  $\eta$  is given by  $\hat{\eta} = M^T\hat{\gamma}$ , where  $\hat{\gamma}$  is the least squares estimate of  $\gamma$ . If  $\epsilon \sim N_n(\mathbf{0}, \sigma^2 I)$ , then  $\hat{\eta} \sim N_s(\eta, \sigma^2 M^T \text{Cov}[\hat{\gamma}] M)$ , where  $\text{Cov}[\hat{\gamma}]$  is as in [Lemma 11.7.2](#) and it can be estimated by  $s^2(\hat{\eta}) = MSE \left[ M^T (\tilde{Z}^T \tilde{Z})^{-1} M \right]$ . We can now easily handle all the inferential issues regarding  $\eta$  such as construction of confidence intervals and tests of hypotheses. In particular, suppose we want to test  $H_0: \eta = \eta_0$  vs  $H_1: \eta \neq \eta_0$ , where  $\eta_0$  is known. A test statistic is

$$F = (\hat{\eta} - \eta_0)^T [s^2(\hat{\eta})]^{-1} (\hat{\eta} - \eta_0)/s.$$

Degrees of freedom associated with this  $F$ -test are  $(s, n - p - q)$ .

#### *Inference for $\beta$*

If we want to estimate  $\theta = L^T\beta$ , where  $L$  is a  $p \times r$  matrix of rank  $r \leq p$ , then its least squares estimate is given by  $\hat{\theta} = L^T\hat{\beta}$  and the distribution of  $\hat{\theta}$  is an  $r$ -dim normal with mean  $\theta$  and covariance matrix  $L^T \text{Cov}[\hat{\beta}] L$ , where the expression for  $\text{Cov}[\hat{\beta}]$  is given in [Lemma 11.7.2](#).

An estimate of  $\text{Cov}[\hat{\theta}]$  is

$$s^2(\hat{\theta}) = MSE L^T \left[ (X^T X)^{-1} + (X^T X)^{-1} X^T D X (X^T X)^{-1} \right].$$

In order to carry out a test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ ,  $\theta_0$  known, we can use the following test statistic

$$F = (\hat{\theta} - \theta_0)^T [s^2(\hat{\theta})]^{-1} (\hat{\theta} - \theta_0) / r.$$

The degrees of freedom for this  $F$ -test are  $(r, n - p - q)$ .

We should point that we can also carry out tests for  $\theta$  or  $\eta$  using the alternate expression for the  $F$ -statistic which involves obtaining  $SSE_R$  and  $SSE_F$ , the residual sums of squares for the reduced and full models, as described in [Section 11.6](#).

**Example 11.7.1.** Consider the following model with one factor and one covariate

$$Y_{ij} = \mu + \alpha_i + \gamma Z_{ij} + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k.$$

In [Example 11.4.6](#), estimates of  $\mu$ ,  $\alpha_i$ , and  $\gamma$  are given. Here  $Z$  has only one column,  $Q_X Z$  is the vector of fitted values when fitting the model  $Z_{ij} = \mu + \alpha_i + \text{error}$ , and  $\tilde{Z} = (I - Q_X)Z$  is the vector of residuals. Hence  $\tilde{Z}_{ij} = Z_{ij} - \bar{Z}_{..}$ , where  $\bar{Z}_{..} = \sum_{j=1}^{n_i} Z_{ij}/n_i$ . Similarly, let  $\tilde{Y}_{ij} = Y_{ij} - \bar{Y}_{..}$ , where  $\bar{Y}_{..} = \sum_{j=1}^{n_i} Y_{ij}/n_i$ . Consequently,

$$\begin{aligned} \hat{\gamma} &= (\tilde{Z}^T \tilde{Z})^{-1} \tilde{Z}^T Y = \left( \sum_i \sum_j \tilde{Z}_{ij} Y_{ij} \right) / \left( \sum_i \sum_j \tilde{Z}_{ij}^2 \right) \\ &= \left( \sum_i \sum_j \tilde{Z}_{ij} \tilde{Y}_{ij} \right) / \left( \sum_i \sum_j \tilde{Z}_{ij}^2 \right). \end{aligned}$$

The residual sum of squares is

$$\begin{aligned} SSE &= \tilde{Y}^T \tilde{Y} - (\tilde{Z}^T \tilde{Z})(\tilde{Z}^T \tilde{Z})^{-1} (\tilde{Z}^T \tilde{Y}) \\ &= \sum_i \sum_j \tilde{Y}_{ij}^2 - \left( \sum_i \sum_j \tilde{Z}_{ij} \tilde{Y}_{ij} \right)^2 / \left( \sum_i \sum_j \tilde{Z}_{ij}^2 \right). \end{aligned}$$

Note that  $df(SSE) = n - k - 1$ , where  $n = \sum n_i$  is the total number of observations. So an estimate of  $\sigma^2$  is  $MSE = SSE/(n - k - 1)$ .

If we want to test  $H_0: \alpha_1 = \dots = \alpha_k = 0$  vs  $H_1: \text{not all } \alpha_i \text{'s are equal to zero}$ , then we need to obtain the residual sum of squares  $SSE_R$  of the reduced model.

The reduced model (under  $H_0$ ) is  $Y_{ij} = \mu + \gamma Z_{ij} + \varepsilon_{ij}$  and

$$\begin{aligned} SSE_R &= \inf_{\mu, \gamma} \sum_i \sum_j [Y_{ij} - \mu - \gamma Z_{ij}]^2 \\ &= \sum_i \sum_j [(Y_{ij} - \bar{Y}_{..}) - \gamma^*(Z_{ij} - \bar{Z}_{..})]^2, \text{ where} \end{aligned}$$

$$\gamma^* = \sum_i \sum_j (Z_{ij} - \bar{Z}_{..}) Y_{ij} \Bigg/ \sum_i \sum_j (Z_{ij} - \bar{Z}_{..})^2,$$

$$df(SSE_R) = n - 2.$$

So the test statistic is

$$F = \frac{(SSE_R - SSE_F)/(k-1)}{MSE},$$

where  $MSE = MSE_F$  and the degrees of freedom for this test are  $(k-1, n-k-1)$ .

If we wish to construct a confidence interval for  $\theta = \alpha_i - \alpha_{i'}$ ,  $i \neq i'$ , then its estimate is  $\hat{\theta} = \bar{Y}_{i.} - \bar{Y}_{i'.} - \hat{\gamma}(\bar{Z}_{i.} - \bar{Z}_{i'.})$ . It is left to the reader to verify that

$$\text{Cov}[\bar{Y}_{i.}, \hat{\gamma}] = 0 \quad \text{for any } i = 1, \dots, k.$$

Fairly simple calculations will show that

$$\begin{aligned} E[\hat{\theta}] &= \theta \text{ and} \\ \text{Var}[\hat{\theta}] &= \sigma^2 \left[ 1/n_i + 1/n_{i'} + (\bar{Z}_{i.} - \bar{Z}_{i'.})^2 / S_{\bar{Z}} \right], \text{ where} \\ S_{\bar{Z}} &= \sum_i \sum_j \tilde{Z}_{ij}^2. \end{aligned}$$

Thus an estimate of  $\text{Var}[\hat{\theta}]$  is given by

$$s^2(\hat{\theta}) = MSE \left[ 1/n_i + 1/n_{i'} + (\bar{Z}_{i.} - \bar{Z}_{i'.})^2 / S_{\bar{Z}} \right].$$

Since  $(\hat{\theta} - \theta)/s(\hat{\theta}) \sim t_{n-k-1}$ , we can construct a confidence interval for  $\theta$ . As a matter of fact, we can carry out many pairwise comparisons  $\alpha_i - \alpha_{i'}$ ,  $i \neq i'$ , using the Bonferroni or Scheffé methods.

#### 11.7.4 Application of ANCOVA to Missing Data

Suppose that we have the usual linear model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{Y}$  is  $n \times 1$  and  $\mathbf{X}$  is  $n \times p$ . However, the last  $r$  observations on the response (ie,  $Y$  values) are missing. Let  $\mathbf{Y}_*$  and  $\mathbf{Y}_r$  denote the vectors of  $n-r$  available observations and the vectors of  $r$  missing observations, respectively. Similarly, let  $\mathbf{X}_*$  and  $\mathbf{X}_r$  denote the submatrices of  $\mathbf{X}$  corresponding to the available and missing cases. Note that  $\mathbf{X}_*$  is  $(n-r) \times p$  and  $\mathbf{X}_r$  is  $r \times p$ . The basic idea

here is to formulate an ANCOVA by introducing one unknown parameter for every missing observation. Consider the ANCOVA model

$$\begin{pmatrix} \mathbf{Y}_* \\ \mathbf{Y}_r \end{pmatrix} = \begin{pmatrix} \mathbf{X}_* \\ \mathbf{X}_r \end{pmatrix} \boldsymbol{\beta} + \begin{pmatrix} \mathbf{0} \\ \mathbf{I}_r \end{pmatrix} \boldsymbol{\gamma} + \boldsymbol{\varepsilon},$$

where  $\mathbf{I}_r$  is the  $r \times r$  identity matrix, and the  $r \times 1$  vector of unknown parameters  $\boldsymbol{\gamma}$  is introduced here for the missing cases. If we carry out a straightforward least squares with the design matrix  $\begin{pmatrix} \mathbf{X}_* & \mathbf{0} \\ \mathbf{X}_r & \mathbf{I}_r \end{pmatrix}$  in order to estimate the unknown parameters  $\begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix}$ , then the normal equations are

$$\begin{pmatrix} \mathbf{X}_*^T \mathbf{X}_* + \mathbf{X}_r^T \mathbf{X}_r & \mathbf{X}_r^T \\ \mathbf{X}_r & \mathbf{I}_r \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}_*^T \mathbf{Y}_* + \mathbf{X}_r^T \mathbf{Y}_r \\ \mathbf{Y}_r \end{pmatrix}.$$

The two sets of equations (obtainable from the normal equations above) are

$$\begin{aligned} \mathbf{X}_*^T \mathbf{X}_* \hat{\boldsymbol{\beta}} + \mathbf{X}_r^T \mathbf{X}_r \hat{\boldsymbol{\beta}} + \mathbf{X}_r^T \hat{\boldsymbol{\gamma}} &= \mathbf{X}_*^T \mathbf{Y}_* + \mathbf{X}_r^T \mathbf{Y}_r, \\ \mathbf{X}_r \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\gamma}} &= \mathbf{Y}_r. \end{aligned}$$

Premultiply the second equation above by  $\mathbf{X}_r^T$  and then subtract it from the first equation to obtain

$$\begin{aligned} \mathbf{X}_*^T \mathbf{X}_* \hat{\boldsymbol{\beta}} &= \mathbf{X}_*^T \mathbf{Y}_*, \text{ and hence} \\ \hat{\boldsymbol{\gamma}} &= \mathbf{Y}_r - \mathbf{X}_r \hat{\boldsymbol{\beta}}. \end{aligned}$$

We now have explicit expressions for  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$ . The vectors of fitted values and residuals are

$$\begin{aligned} \hat{\mathbf{Y}} &= \begin{pmatrix} \mathbf{X}_* \\ \mathbf{X}_r \end{pmatrix} \hat{\boldsymbol{\beta}} + \begin{pmatrix} \mathbf{0} \\ \mathbf{I}_r \end{pmatrix} \hat{\boldsymbol{\gamma}} = \begin{pmatrix} \mathbf{X}_* \hat{\boldsymbol{\beta}} \\ \mathbf{X}_r \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\gamma}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}_* \hat{\boldsymbol{\beta}} \\ \mathbf{Y}_r \end{pmatrix}, \\ \mathbf{Y} - \hat{\mathbf{Y}} &= \begin{pmatrix} \mathbf{Y}_* \\ \mathbf{Y}_r \end{pmatrix} - \begin{pmatrix} \mathbf{X}_* \hat{\boldsymbol{\beta}} \\ \mathbf{Y}_r \end{pmatrix} = \begin{pmatrix} \mathbf{Y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}} \\ \mathbf{0} \end{pmatrix}. \end{aligned}$$

Consequently, the residual sum of squares is given by

$$SSE = \|\mathbf{Y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}\|^2.$$

Note that this residual sum of squares turns out to be identical to that for the model ( $\mathbf{Y}_* = \mathbf{X}_* \boldsymbol{\beta} + \boldsymbol{\varepsilon}_*$ ) with only available observations. The degrees of freedom for the SSE is  $n-p-r$ . Thus, the degrees of freedom of the SSE is the same as the one that can be obtained from an analysis of the available observations assuming that  $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}_*)$ .

## 11.8 Model Selection

In the Gauss-Markov setup one is often concerned with selecting a model from an appropriate class of candidate models. There are a number of methods for model selection and properties of these methods have been investigated by many authors, but we focus on a few of them instead of describing all. Consider the general framework  $\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\varepsilon}$ , where  $\mathbf{Y}$

is the  $n \times 1$  response vector,  $\mu$  is the mean vector, and  $\epsilon$  is the vector of iid errors with mean 0 and variance  $\sigma^2$ . We assume that there is a class of candidate models  $\{\mathbf{X}_k \beta_k: k = 1, \dots, K\}$  for describing  $\mu$ , where  $\mathbf{X}_k$  is a known  $n \times p_k$  matrix of rank  $p_k$  and  $\beta_k$  is unknown.

### 11.8.1 An Overview of Various Model Selection Criteria

For the  $k$ th model, if  $\hat{\mu}_k = \mathbf{X}_k \hat{\beta}_k$  is the estimated mean vector, where  $\hat{\beta}_k$  is the least squares estimate of  $\beta_k$ , then the residual sum of squares is  $SSE_k = \|\mathbf{Y} - \hat{\mu}_k\|^2$ . Note that the residual sum of squares  $SSE_k = \|\mathbf{Y} - \mathbf{X}_k \hat{\beta}_k\|^2$  decreases for a nested class of models (where for each  $k$ , the  $k$ th model is contained in the  $k+1$ -st model) and the minimum of  $\{SSE_k\}$  is attained for the model with the largest number of parameters. Thus one cannot hope to select an appropriate model by minimizing  $SSE_k$  over  $k$ . A reasonable measure of how well  $\hat{\mu}_k$  estimates  $\mu$  is given by  $D_k = E[\|\mu - \hat{\mu}_k\|^2]$ , the expected squared distance between  $\mu$  and  $\hat{\mu}_k$ . Ideally, one would choose the model for which  $D_k$  is the smallest, but  $D_k$  is unknown as it depends on the unknown population parameters. Therefore, one then first obtains a good estimate  $\hat{D}_k$  of  $D_k$ , then minimizes  $\hat{D}_k$  over  $k = 1, \dots, K$ . If the minimum is attained at  $k = \hat{k}$ , then  $\{\mathbf{X}_{\hat{k}} \hat{\beta}_{\hat{k}}\}$  is considered the most appropriate estimate of the mean vector  $\mu$ . Akaike's FPE and Mallows' criteria seek to obtain an unbiased (or nearly unbiased) estimate of  $D_k$ .

Cross-validation (CV) and generalized cross-validation (GCV) seek to estimate the prediction error of the  $k$ th fitted model. In this setup, it is assumed that  $(Y_i, \mathbf{x}_{k,i})$ ,  $i = 1, \dots, n$ , where the observed values of the covariates  $\{\mathbf{x}_{k,i}^T\}$  are the rows of  $\mathbf{X}_k$ , are iid and  $(Y_{n+1}, \mathbf{x}_{k,n+1})$  is an independent copy of  $(Y_i, \mathbf{x}_{k,i})$ . If  $\hat{\beta}_k$  is the least squares estimate of  $\beta_k$ , then the error for predicting  $Y_{n+1}$  using the estimated  $k$ th model is  $PE(k) = E[Y_{n+1} - \hat{\beta}_k^T \mathbf{x}_{k,n+1}]^2$ . Both CV and GCV obtain nearly unbiased estimators of this prediction error. Thus one chooses the model for which  $\widehat{PE}(k)$  (as given by the CV or GCV criterion) is the smallest.

Both Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) are likelihood-based methods, which seek to model the joint pdf  $f$  of  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ . Suppose there is a class of candidate probability models  $\{f_k(\cdot, \theta_k): k = 1, \dots, K\}$  for  $f$  under consideration. In the case of linear models,  $f_k$  may be the pdf of  $N_n(\mathbf{X}_k \beta_k, \sigma^2 \mathbf{I})$  and  $\theta_k = \begin{pmatrix} \beta_k \\ \sigma_k \end{pmatrix}$ . As in the case of residual sum of squares,  $-2 \log f(\mathbf{Y}, \hat{\theta}_k)$ , where  $\hat{\theta}_k$  is the MLE for  $\theta_k$ , will always decrease if we take a sequence of nested models and thus it is not possible to select an appropriate model by trying to minimize  $-2 \log f(\mathbf{Y}, \hat{\theta}_k)$  over  $k$ , since the minimum is always attained at the model with the largest number of parameters. A true measure of how good the  $k$ th fitted model is can be judged by how well it performs for a dataset that is independent of the data  $\mathbf{Y}$  and this is what is done in the arguments involving the AIC. The AIC seeks to obtain an estimate of 2 times the Kullback-Leibler

divergence  $KL(k) = E[\log f(\tilde{Y})/f_k(\tilde{Y}, \hat{\theta}_k)]$ , where  $\tilde{Y}$  has the same distribution as  $Y$ , but is independent of it. The model corresponding to the smallest value of AIC is considered the most appropriate one.

As the name suggests, the BIC is based on Bayesian considerations. In this setting, there is an “index of model” variable  $J$  which takes values in  $\{1, \dots, K\}$  so that the prior probability of choosing model  $k$  is  $\pi_k$  (ie,  $P[J = k] = \pi_k$ ). Given  $J = k$ , the pdf of  $Y$  is  $f_k(\cdot, \theta_k)$ , where  $\theta_k$  has a prior  $g_k$ . The goal to find the model for which the conditional probability  $P[J = k|Y]$  is maximized or  $-2 \log P[J = k|Y]$  is minimized. Since these conditional probabilities are unknown, the BIC obtains an estimate of  $-2 \log P[J = k|Y]$  for each  $k$ .

For the  $k$ th linear model, let  $SSE_k = \|Y - X_k \hat{\beta}_k\|^2$  and  $\hat{\sigma}_k^2 = SSE_k/n$ , where  $\hat{\beta}_k$  is the least squares estimate of  $\beta_k$ . Under normality (ie,  $\epsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$ ), the MLE of  $\sigma^2$  is  $\hat{\sigma}_k^2$ . The various model selection criteria mentioned above are

$$\begin{aligned} FPE(k) &= \frac{n + p_k}{n - p_k} SSE_k, \\ MAL(k) &= SSE_k + 2p_k \hat{\sigma}^2, \\ AIC(k) &= n \log \hat{\sigma}_k^2 + 2p_k, \\ BIC(k) &= n \log \hat{\sigma}_k^2 + \log(n)p_k, \\ \widehat{PE}^{(cv)}(k) &= n^{-1} \sum_{i=1}^n \left( Y_i - \hat{\beta}_{k,(-i)}^T \mathbf{x}_{k,i} \right)^2, \text{ and} \\ \widehat{PE}^{(gcv)}(k) &= n^{-1} SSE_k / (1 - p_k/n)^2. \end{aligned}$$

For Mallows’ criterion, it is assumed that  $\hat{\sigma}^2$  is a consistent estimate of  $\sigma^2$ . In the cross-validation criterion, it is understood that, for the  $k$ th model, the rows of  $X_k$  are  $\{\mathbf{x}_{k,i}^T, i = 1, \dots, n\}$  and  $\hat{\beta}_{k,(-i)}$  is the least squares estimate of  $\beta_k$  based on  $n - 1$  observations, deleting  $(Y_i, \mathbf{x}_{k,i})$ .

*Remark 11.8.1.* Both AIC and BIC have forms which are more general than what are written above. More general versions are given below and they are described in detail. It can be shown that FPE, Mallows’, CV, GCV, and AIC criteria are equivalent in an asymptotic sense as  $n \rightarrow \infty$ . The BIC is different from the others as its use may lead to models with fewer parameters. If the correct model is in the candidate class, then mathematical arguments, under appropriate conditions, show that BIC selects the correct model with probability converging to 1. The other model selection criteria given above tend to choose the “best” predictive model. It is important to keep in mind that the correct model (if it exists) is not necessarily the best predictive model.

### 11.8.2 Akaike’s FPE and Mallows’ Criteria

Suppose  $\hat{\beta}_k$  is the least squares estimate of  $\beta_k$  and  $\hat{\mu}_k = X_k \hat{\beta}_k$ . A measure of divergence between  $\mu$  and  $\hat{\mu}_k$  is given by  $D_k = E[\|\mu - \hat{\mu}_k\|^2]$ , and the goal is to obtain a reasonable

estimate of  $D_k$ . Now we can write  $\hat{\mu}_k = \mathbf{Q}_k \mathbf{Y}$ , where  $\mathbf{Q}_k = \mathbf{X}_k (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{X}_k^T$  is the projection on the column space of  $\mathbf{X}_k$  and  $\text{rank}(\mathbf{Q}_k) = \text{trace}(\mathbf{Q}_k) = p_k$ . Moreover,

$$\boldsymbol{\mu} - \hat{\mu}_k = \boldsymbol{\mu} - \mathbf{Q}_k \mathbf{Y} = (\mathbf{I} - \mathbf{Q}_k) \boldsymbol{\mu} - \mathbf{Q}_k \boldsymbol{\varepsilon} = \bar{\mathbf{Q}}_k \boldsymbol{\mu} - \mathbf{Q}_k \boldsymbol{\varepsilon},$$

where  $\bar{\mathbf{Q}}_k = \mathbf{I} - \mathbf{Q}_k$ . Since  $\bar{\mathbf{Q}}_k \boldsymbol{\mu}$  and  $\mathbf{Q}_k \boldsymbol{\varepsilon}$  are orthogonal,  $\|\boldsymbol{\mu} - \hat{\mu}_k\|^2 = \|\bar{\mathbf{Q}}_k \boldsymbol{\mu}\|^2 + \|\mathbf{Q}_k \boldsymbol{\varepsilon}\|^2$  and taking expectation on both sides, we have

$$D_k = E\|\boldsymbol{\mu} - \hat{\mu}_k\|^2 = \|\bar{\mathbf{Q}}_k \boldsymbol{\mu}\|^2 + p_k \sigma^2.$$

Let us now examine the residual sum of squares  $SSE_k = \|\mathbf{Y} - \hat{\mu}_k\|^2$ . We can write  $\mathbf{Y} - \hat{\mu}_k = \bar{\mathbf{Q}}_k \mathbf{Y} = \bar{\mathbf{Q}}_k \boldsymbol{\mu} + \bar{\mathbf{Q}}_k \boldsymbol{\varepsilon}$ . Hence

$$\begin{aligned} SSE_k &= \|\bar{\mathbf{Q}}_k \boldsymbol{\mu}\|^2 + \|\bar{\mathbf{Q}}_k \boldsymbol{\varepsilon}\|^2 + 2(\bar{\mathbf{Q}}_k \boldsymbol{\mu})^T (\bar{\mathbf{Q}}_k \boldsymbol{\varepsilon}), \text{ and} \\ E[SSE_k] &= \|\bar{\mathbf{Q}}_k \boldsymbol{\mu}\|^2 + (n - p_k) \sigma^2. \end{aligned}$$

This shows that

$$D_k = E[SSE_k] - (n - 2p_k) \sigma^2 = E[SSE_k] + 2p_k \sigma^2 - n \sigma^2.$$

The last term  $n \sigma^2$  does not depend on  $k$  and hence there is no need to estimate it. So if we have an estimate  $\hat{\sigma}_k^2$  of  $\sigma^2$ , then we can get a criterion

$$\hat{D}_k = SSE_k + 2p_k \hat{\sigma}_k^2 - n \sigma^2.$$

The criteria given in [Section 11.8.1](#) ignore the  $n \sigma^2$  term. How one estimates  $\sigma^2$  depends on the type of problem at hand. If  $\sigma^2$  is estimated by  $MSE_k = SSE_k / (n - p_k)$ , then one ends up with Akaike's FPE criterion given above. Another possibility is to take a large enough model (maybe the largest model in the candidate class) so that the model bias is low and use the mean square error of that model to estimate  $\sigma^2$ . In such a case, the estimate of  $\sigma^2$  does not depend on  $k$  and the resulting  $\hat{D}_k$  is a special case of Mallows' criterion.

### 11.8.3 AIC and BIC

Let  $f$  be the joint pdf of  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  and the goal is to find a suitable approximation for  $f$  among a class of candidate probability models. Let the  $k$ th model be  $\{f_k(y, \theta_k), \theta_k \in \Omega_k\}$ ,  $k = 1, \dots, K$ , where  $\Omega_k$  is a nonempty open subset of  $\mathbb{R}^{p_k}$ . If  $\hat{\theta}_k$  is the MLE of  $\theta_k$  for the  $k$ th model, then the AIC and BIC are

$$\begin{aligned} AIC(k) &= -2 \log f_k(\mathbf{Y}, \hat{\theta}_k) + 2p_k, \\ BIC(k) &= -2 \log f_k(\mathbf{Y}, \hat{\theta}_k) + \log(n)p_k. \end{aligned}$$

If a criterion is minimized at  $k = \hat{k}$ , then  $f_{\hat{k}}$  is declared to be the most appropriate model according to that criterion. Note that in both the criteria,  $-2 \log f_k(\mathbf{Y}, \hat{\theta}_k)$  are penalized by a constant times the number of parameters estimated. However, the BIC has a higher penalty than the AIC and thus it may select a model with fewer parameters than the AIC.

In the case of linear models with normal errors, if one has a sequence of models  $\{\mathbf{X}_k \boldsymbol{\beta}_k : \boldsymbol{\beta}_k \in \mathbb{R}^{p_k}\}$ , then there are simple expressions for these criteria. Let  $\boldsymbol{\theta}_k = \begin{pmatrix} \boldsymbol{\beta}_k \\ \sigma_k \end{pmatrix}$  be the vector of parameters for the  $k$ th model, then the MLE are

$$\hat{\boldsymbol{\beta}}_k = (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{X}_k^T \mathbf{Y}, \quad \hat{\sigma}_k^2 = \|\mathbf{Y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k\|^2/n,$$

and thus we have

$$\begin{aligned} f_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) &= \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}_k} \right)^n \exp \left[ -\|\mathbf{Y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k\|^2 / (2\hat{\sigma}_k^2) \right] \\ &= \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}_k} \right)^n \exp[-n/2], \text{ and} \\ -2 \log f_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) &= n \log \hat{\sigma}_k^2 + n + n \log(2\pi). \end{aligned}$$

The last two terms in the last expression do not involve  $k$ , and after ignoring them we arrive at

$$AIC(k) = n \log \hat{\sigma}_k^2 + 2p_k \quad \text{and} \quad BIC(k) = n \log \hat{\sigma}_k^2 + \log(n)p_k.$$

Theoretical arguments used in deriving these criteria are different. The arguments given here are heuristic, but they can be made rigorous at the expense of fairly cumbersome calculations. The AIC is obtained by trying to approximate a predictive likelihood, whereas the derivation of the BIC involves approximating the probability of choosing a model given the data  $\mathbf{Y}$ .

### *Mathematical Settings*

- I. Suppose  $\tilde{\mathbf{Y}}$  has the same distribution as  $\mathbf{Y}$  but is independent of it. AIC is obtained by an approximation of  $E[-2 \log f(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)]$ , where  $\hat{\boldsymbol{\theta}}_k$  is the MLE based on the available data.
- II. Suppose  $J$  is a discrete rv taking values  $1, \dots, K$ , and  $P[J = k] = \pi_k$ . Also assume that given  $J = k$ , the pdf of the model is  $f(\cdot, \boldsymbol{\theta}_k)$ , where  $\boldsymbol{\theta}_k$  has a prior  $g_k(\cdot)$ . Then the goal is to find the value of  $k$  which maximizes  $P[J = k | \mathbf{Y}]$  or minimizes  $-2 \log P[J = k | \mathbf{Y}]$ . The BIC is obtained by finding an approximation to  $-2 \log P[J = k | \mathbf{Y}]$ . It turns out that as long as the priors  $\{\pi_k\}$  and  $\{g_k(\cdot)\}$  are reasonably “flat” on the parameter spaces (ie, they are not degenerate or close to being degenerate), the dominant terms in the asymptotic expansion of  $-2 \log P[J = k | \mathbf{Y}]$  do not depend on the priors.

### *Heuristic Derivation of AIC*

Let us assume that there exists a unique  $\boldsymbol{\theta}_{k0} \in \Omega_k$  when  $E[\log f(\mathbf{Y}, \boldsymbol{\theta}_k)]$  is maximized over  $\boldsymbol{\theta}_k \in \Omega_k$ , where the expectation is over the true pdf  $f$  of  $\mathbf{Y}$ . Let  $f_i$  be the pdf of  $Y_i$  and  $f_{k,i}(\cdot, \boldsymbol{\theta}_k)$  be the pdf of  $Y_i$  under the  $k$ th model and denote

$$\begin{aligned}\xi_{k,i}(\cdot, \boldsymbol{\theta}_k) &= -\log f_{k,i}(\cdot, \boldsymbol{\theta}_k), \quad \xi_k(\mathbf{y}, \boldsymbol{\theta}_k) = \sum_{i=1}^n \xi_{k,i}(y_i, \boldsymbol{\theta}_k), \\ \dot{\xi}_{k,i}(\cdot, \boldsymbol{\theta}_k) &= \frac{\partial}{\partial \boldsymbol{\theta}_k} \xi_{k,i}(\cdot, \boldsymbol{\theta}_k), \quad \ddot{\xi}_{k,i}(\cdot, \boldsymbol{\theta}_k) = \frac{\partial^2}{\partial \boldsymbol{\theta}_k \partial \boldsymbol{\theta}_k} \xi_{k,i}(\cdot, \boldsymbol{\theta}_k), \\ \dot{\xi}_k(\mathbf{y}, \boldsymbol{\theta}_k) &= \sum_{i=1}^n \dot{\xi}_{k,i}(y_i, \boldsymbol{\theta}_k), \quad \text{and} \quad \ddot{\xi}_k(\mathbf{y}, \boldsymbol{\theta}_k) = \sum_{i=1}^n \ddot{\xi}_{k,i}(y_i, \boldsymbol{\theta}_k).\end{aligned}$$

In these notations,  $f(\mathbf{y}, \boldsymbol{\theta}_k) = \exp[-\xi_k(\mathbf{y}, \boldsymbol{\theta}_k)]$ . Let

$$\mathbf{I}_k(f) = \mathbb{E} \left[ n^{-1} \sum_{i=1}^n \dot{\xi}_{k,i}(Y_i, \boldsymbol{\theta}_{k0}) \dot{\xi}_{k,i}(Y_i, \boldsymbol{\theta}_{k0})^T \right], \quad \tilde{\mathbf{I}}_k(f) = \mathbb{E} \left[ n^{-1} \ddot{\xi}_k(Y, \boldsymbol{\theta}_{k0}) \right],$$

where the expectations are taken over the true pdf  $f(\cdot)$  of  $\mathbf{Y}$ . Suppose that  $\hat{\boldsymbol{\theta}}_k$  is the MLE of  $\boldsymbol{\theta}_k$  based on the data  $\mathbf{Y}$ ,  $\tilde{\mathbf{Y}}$  is an independent copy of  $\mathbf{Y}$ , and the following hold (heuristic justifications given below)

$$\mathbb{E}[\xi_k(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)] = \mathbb{E}[\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k)] + (1/2) \text{trace}(\tilde{\mathbf{I}}_k(f)^{-1} \mathbf{I}_k(f)) [1 + o(1)], \quad (5a)$$

$$\mathbb{E}[\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k)] = \mathbb{E}[\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k)] - (1/2) \text{trace}(\tilde{\mathbf{I}}_k(f)^{-1} \mathbf{I}_k(f)) [1 + o(1)]. \quad (5b)$$

Then a reasonable estimate of  $\mathbb{E}[\xi_k(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)]$  is given by

$$\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) + \text{trace}(\tilde{\mathbf{I}}_k(f)^{-1} \mathbf{I}_k(f)).$$

Since  $\text{trace}(\tilde{\mathbf{I}}_k(f)^{-1} \mathbf{I}_k(f))$  is unknown, it needs to be estimated. If it is assumed that  $f_k(\cdot, \boldsymbol{\theta}_{k0})$  is reasonably close to  $f$ , then we may replace  $\tilde{\mathbf{I}}_k(f)$  and  $\mathbf{I}_k(f)$  by  $\tilde{\mathbf{I}}_k(f_k(\cdot, \boldsymbol{\theta}_{k0}))$  and  $\mathbf{I}(f_k(\cdot, \boldsymbol{\theta}_{k0}))$ , respectively. Nothing that  $\tilde{\mathbf{I}}_k(f_k(\cdot, \boldsymbol{\theta}_{k0})) = \mathbf{I}_k(f_k(\cdot, \boldsymbol{\theta}_{k0}))$  (using Eqs. (3a) and (3b) in Section 7.1, and  $\xi_k$  is the negative of the log-likelihood), we may estimate  $\text{trace}(\tilde{\mathbf{I}}_k(f)^{-1} \mathbf{I}_k(f))$  by

$$\text{trace}(\tilde{\mathbf{I}}_k(f_k(\cdot, \boldsymbol{\theta}_{k0}))^{-1} \mathbf{I}_k(f_k(\cdot, \boldsymbol{\theta}_{k0}))) = \text{trace}(\mathbf{I}) = p_k.$$

Thus an estimate of  $\mathbb{E}[\xi_k(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)]$  is given by

$$\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) + p_k = -\log f_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) + p_k.$$

The AIC is two times the quantity given above.

What remains to be shown is that the approximate expansions given in Eqs. (5a) and (5b) are valid. The MLE  $\hat{\boldsymbol{\theta}}_k$  satisfies the likelihood equation  $\dot{\xi}_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) = \mathbf{0}$ . Since  $\boldsymbol{\theta}_{k0}$  maximizes  $\mathbb{E}[\log f(\mathbf{Y}, \boldsymbol{\theta}_k)]$ ,  $\mathbb{E}[\dot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})] = 0$ . For notational convenience, we denote  $\mathbf{I}_k(f)$

and  $\tilde{\mathbf{I}}_k(f)$  by  $\mathbf{I}_k$  and  $\tilde{\mathbf{I}}_k$ . Arguments used in obtaining the asymptotic properties of the MLE (Chapter 7) can be employed to conclude the following

$$\begin{aligned} n^{-1}\ddot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0}) - \tilde{\mathbf{I}}_k &\xrightarrow{P} \mathbf{0}, \\ \mathbf{Z}_{n,k} = n^{-1/2}\mathbf{I}_k^{-1/2}\dot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0}) &\xrightarrow{\mathcal{D}} \mathbf{Z}_k \sim N_{p_k}(\mathbf{0}, \mathbf{I}), \text{ and} \\ \hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0} &= -\ddot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})^{-1}\dot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})[1 + o_P(1)] \\ &= -n^{-1/2}\tilde{\mathbf{I}}_k^{-1}\mathbf{I}_k^{1/2}\mathbf{Z}_{n,k}[1 + o_P(1)]. \end{aligned}$$

Now

$$\begin{aligned} \xi_k(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k) &= \xi_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0}) + (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\dot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0}) \\ &\quad + (1/2)(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\ddot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})[1 + o_P(1)]. \end{aligned}$$

Since  $\tilde{\mathbf{Y}}$  is independent of  $\mathbf{Y}$  and hence of  $\hat{\boldsymbol{\theta}}_k$ ,

$$\mathbb{E}[\xi_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})] = \mathbb{E}[\xi_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})], \quad \mathbb{E}[\dot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})] = \mathbb{E}[\dot{\xi}(\mathbf{Y}, \boldsymbol{\theta}_{k0})] = \mathbf{0}.$$

Noting that  $n^{-1}\ddot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0}) - \tilde{\mathbf{I}}_k \xrightarrow{P} \mathbf{0}$ , we have

$$\begin{aligned} &(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\ddot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0}) \\ &= \mathbf{Z}_{n,k}^T \mathbf{I}_k^{1/2} \tilde{\mathbf{I}}_k^{-1} [n^{-1}\ddot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})] \tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k^{1/2} \mathbf{Z}_{n,k} [1 + o_P(1)] \\ &= \mathbf{Z}_{n,k}^T \mathbf{I}_k^{1/2} \tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k^{1/2} \mathbf{Z}_{n,k} [1 + o_P(1)]. \end{aligned}$$

Since  $\mathbf{Z}_{n,k} \xrightarrow{\mathcal{D}} \mathbf{Z}_k$  and  $\mathbb{E}[\mathbf{Z}_k \mathbf{Z}_k^T] = \mathbf{I}$ , ignoring the  $o_P(1)$  term we have

$$\begin{aligned} \mathbb{E}\left[(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\ddot{\xi}_k(\tilde{\mathbf{Y}}, \boldsymbol{\theta}_{k0})(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})\right] &= \mathbb{E}\left[\mathbf{Z}_k^T \mathbf{I}_k^{1/2} \tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k^{1/2} \mathbf{Z}_k\right] [1 + o(1)] \\ &= \text{trace}(\mathbf{I}_k^{1/2} \tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k^{1/2} \mathbb{E}[\mathbf{Z}_k \mathbf{Z}_k^T]) [1 + o(1)] \\ &= \text{trace}(\tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k) [1 + o(1)], \text{ and} \\ \mathbb{E}[\xi_k(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)] &= \mathbb{E}[\xi_k(\mathbf{Y}, \boldsymbol{\theta}_k)] + (1/2)\text{trace}(\tilde{\mathbf{I}}_k^{-1} \mathbf{I}_k) [1 + o(1)]. \end{aligned}$$

This shows that Eq. (5a) holds. The above argument assumes that the  $o_P(1)$  term can be turned into an  $o(1)$  term when taking the expectation. Such a step requires the concept of *uniform integrability* which we ignore here for the sake of simplicity of discussion.

In order to verify Eq. (5b), expand  $\xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k)$  about  $\xi_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})$  to get

$$\begin{aligned} \xi_k(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) &= \xi_k(\mathbf{Y}, \boldsymbol{\theta}_{k0}) + (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\dot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0}) \\ &\quad + (1/2)(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})^T\ddot{\xi}_k(\mathbf{Y}, \boldsymbol{\theta}_{k0})(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_{k0})[1 + o_P(1)]. \end{aligned}$$

Since  $\hat{\theta}_k - \theta_{k0} = -\ddot{\xi}_k(\mathbf{Y}, \theta_{k0})^{-1}\dot{\xi}_k(\mathbf{Y}, \theta_{k0})[1 + o_P(1)]$ , we have

$$\xi_k(\mathbf{Y}, \hat{\theta}_k) = \xi_k(\mathbf{Y}, \theta_{k0}) - (1/2)(\hat{\theta}_k - \theta_{k0})^T \ddot{\xi}_k(\mathbf{Y}, \theta_{k0})(\hat{\theta}_k - \theta_{k0})[1 + o_P(1)].$$

Using the same argument as in the expansion of  $\xi_k(\tilde{\mathbf{Y}}, \hat{\theta}_k)$ , we get

$$\mathbb{E}[\xi_k(\mathbf{Y}, \hat{\theta}_k)] = \mathbb{E}[\xi_k(\mathbf{Y}, \theta_{k0})] - (1/2)\text{trace}(\tilde{\mathbf{I}}_k^{-1}\mathbf{I}_k)[1 + o(1)].$$

### *Heuristic Derivation of BIC*

Recall that in the Bayesian setting, the goal is to find  $k$  for which  $P[J = k|\mathbf{Y}]$  is maximized where the discrete variable  $J$  can take values  $1, \dots, K$  with probabilities  $\pi_1, \dots, \pi_K$ . Now

$$P[J = k|\mathbf{Y} = \mathbf{y}] = \frac{\int_{\Omega_k} f_k(\mathbf{y}, \theta_k) g_k(\theta_k) \pi_k d\theta_k}{f(\mathbf{y})},$$

where  $f(\cdot)$  is the marginal pdf of  $\mathbf{Y}$ . Since the denominator of  $f(\mathbf{y})$  does not depend on  $k$ , we may simply ignore it as it plays no role in optimization of  $P[J = k|\mathbf{Y}]$  over  $k$ . We outline the basic arguments in approximating the integral in the numerator when  $g_k$  is taken as a noninformative prior (ie,  $g_k(\theta) \equiv 1$ ). When  $\Omega_k$  is not compact, the use of such a prior can be justified by taking a uniform prior on a compact region  $A_{k,j} \subset \Omega_k$  and letting  $A_{k,j} \rightarrow \Omega_k$  as  $j \rightarrow \infty$  in an appropriate manner.

Let  $\xi_k$ ,  $\xi_{k,i}$ , etc., be the same as in the derivation of the AIC. Then  $f_k(\mathbf{Y}, \theta_k) = \exp[-\xi_k(\mathbf{Y}, \theta_k)]$ . The MLE  $\hat{\theta}_k$  of  $\theta_k$  is a solution to the equation  $\dot{\xi}_k(\mathbf{Y}, \theta_k) = 0$  and assume that  $n^{-1}\ddot{\xi}_k(\mathbf{Y}, \hat{\theta}_k) \xrightarrow{P} \tilde{\mathbf{I}}_k$ , a positive definite matrix. Reexpress  $\theta_k$  as  $\hat{\theta}_k + n^{-1/2}\mathbf{u}$  where  $\mathbf{u} = \sqrt{n}(\theta_k - \hat{\theta}_k)$  and hence

$$\begin{aligned} \int_{\Omega_k} f_k(\mathbf{Y}, \theta_k) d\theta_k &= n^{-p_k/2} \int_{\Omega'_k} f_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u}) d\mathbf{u} \\ &= n^{-p_k/2} \int_{\Omega'_k} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u})] d\mathbf{u}, \end{aligned}$$

where  $\Omega'_k = \left\{ \sqrt{n}(\mathbf{z} - \hat{\theta}_k) : \mathbf{z} \in \Omega_k \right\}$ . If it can be shown that (heuristic details given below)

$$\int_{\Omega'_k} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u})] d\mathbf{u} = \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)] (\sqrt{2\pi})^{p_k} |\tilde{\mathbf{I}}_k|^{-1/2} [1 + o_P(1)],$$

then we have

$$\int_{\Omega_k} f_k(\mathbf{Y}, \theta_k) d\theta_k = n^{-p_k/2} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)] (\sqrt{2\pi})^{p_k} |\tilde{\mathbf{I}}_k|^{-1/2} [1 + o_P(1)].$$

Since  $f_k(\mathbf{Y}, \hat{\theta}_k) = \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)]$ , keeping only the dominant terms,  $-2$  times the logarithm of  $\int_{\Omega'_k} f_k(\mathbf{Y}, \theta_k) d\theta_k$  is approximately equal to

$$2\xi_k(\mathbf{Y}, \hat{\theta}_k) + \log(n)p_k = -2 \log f_k(\mathbf{Y}, \hat{\theta}_k) + \log(n)p_k,$$

which is the BIC.

Let us now outline the heuristic arguments involved in the approximation of the integral  $\int_{\Omega'_k} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u})] d\mathbf{u}$ . Assuming that  $\hat{\theta}_k$  is in a compact set with probability converging to 1,  $\Omega'_k \rightarrow \mathbb{R}^{p_k}$  as  $n \rightarrow \infty$ . Since  $\hat{\theta}_k$  is the MLE, we have  $\dot{\xi}_k(\mathbf{Y}, \hat{\theta}_k) = 0$ . Expanding  $\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u})$  about  $\xi_k(\mathbf{Y}, \hat{\theta}_k)$ , we have (under reasonable regularity conditions on  $f_k$ )

$$\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u}) = \xi_k(\mathbf{Y}, \hat{\theta}_k) + (1/2)\mathbf{u}^T \tilde{\mathbf{I}}_k \mathbf{u}[1 + o_P(1)].$$

Thus

$$\begin{aligned} & \int_{\Omega'_k} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k + n^{-1/2}\mathbf{u})] d\mathbf{u} \\ &= \int_{\Omega'_k} \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k) - (1/2)\mathbf{u}^T \tilde{\mathbf{I}}_k \mathbf{u}(1 + o_P(1))] d\mathbf{u} \\ &= \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)] \int_{\Omega'_k} \exp[-(1/2)\mathbf{u}^T \tilde{\mathbf{I}}_k \mathbf{u}(1 + o_P(1))] d\mathbf{u} \\ &= \exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)] \int_{\mathbb{R}^{p_k}} \exp[-(1/2)\mathbf{u}^T \tilde{\mathbf{I}}_k \mathbf{u}] d\mathbf{u}[1 + o_P(1)]. \end{aligned}$$

This kind of approximation of an integral is known as the Laplace approximation. We have omitted a lot of details in the above approximation including how the integral over  $\Omega'_k$  can be approximated by the integral over  $\mathbb{R}^{p_k}$ , how the  $o_P(1)$  terms can be brought out of the exponential term, etc. However, all the calculations can be justified using appropriate mathematical conditions.

The integrand in the last integral is proportional to the pdf of a  $p_k$ -dim normal random vector with mean  $\mathbf{0}$  and covariance matrix  $\tilde{\mathbf{I}}_k^{-1}$ . Thus

$$\int_{\mathbb{R}^{p_k}} \exp[-(1/2)\mathbf{u}^T \tilde{\mathbf{I}}_k \mathbf{u}] d\mathbf{u} = (\sqrt{2\pi})^{p_k} |\tilde{\mathbf{I}}_k|^{-1/2}.$$

Thus an approximation to  $\int_{\Omega'_k} f_k(\mathbf{Y}, \theta_k) d\theta_k$  is given by

$$\exp[-\xi_k(\mathbf{Y}, \hat{\theta}_k)] (\sqrt{2\pi})^{p_k} |\tilde{\mathbf{I}}_k|^{-1/2}.$$

#### 11.8.4 Cross-Validation and Generalized Cross-Validation

Recall that the  $k$ th model under consideration is  $\mathbf{Y} = \mathbf{X}_k \beta_k + \boldsymbol{\varepsilon}$ , where  $\mathbf{X}_k$  is an  $n \times p_k$  matrix of rank  $p_k$ , and assume that the rows of  $\mathbf{X}_k$  are  $\mathbf{x}_{k,i}^T$ ,  $i = 1, \dots, n$ , and that  $(Y_i, \mathbf{x}_{k,i})$ ,

$i = 1, \dots, n$ , are iid. If  $(Y_{n+1}, \mathbf{x}_{k,n+1})$  is a future copy of  $(Y_i, \mathbf{x}_{k,i})$ , which is independent of the data  $\{(Y_i, \mathbf{x}_{k,i}), i = 1, \dots, n\}$ , then the mean square error in predicting  $Y_{n+1}$  by  $\hat{\beta}_k^T \mathbf{x}_{k,n+1}$  is  $PE(k) = E[Y_{n+1} - \hat{\beta}_k^T \mathbf{x}_{k,n+1}]^2$ , where  $\hat{\beta}_k = (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{X}_k^T \mathbf{Y}$  is the least squares estimate of  $\beta_k$ . Cross-validation method seeks to estimate this prediction error  $PE(k)$ . The method is as follows. Obtain the least squares estimate  $\hat{\beta}_{k,(-i)}$  of  $\beta_k$  based on  $n - 1$  observations deleting the  $i$ th case. Then an estimate of the prediction error is given by

$$\widehat{PE}^{(cv)}(k) = n^{-1} \sum_{i=1}^n \left( Y_i - \hat{\beta}_{k,(-i)}^T \mathbf{x}_{k,i} \right)^2.$$

As will be shown below, there is no need to calculate all the  $n$  different estimates of  $\beta_k$ . Let  $\mathbf{Q}_k = \mathbf{X}_k (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{X}_k^T$  be the projection matrix which projects on the columns of  $\mathbf{X}_k$ . Let  $\hat{Y}_i = \hat{\beta}_k^T \mathbf{x}_i$  and  $\hat{\varepsilon}_i = Y_i - \hat{Y}_i$ . It will be shown that

$$\hat{\beta}_{k,(-i)} = \hat{\beta}_k - (1 - q_{k,ii})^{-1} (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{x}_{k,i} \hat{\varepsilon}_i, \quad (6a)$$

$$Y_i - \hat{\beta}_{k,(-i)}^T \mathbf{x}_{k,i} = \left( Y_i - \hat{\beta}_k^T \mathbf{x}_{k,i} \right) / (1 - q_{k,ii}), \quad (6b)$$

where  $q_{k,ii}$  is the  $i$ th diagonal element of  $\mathbf{Q}_k$ . So  $PE^{(cv)}$  can be written as

$$\widehat{PE}^{(cv)}(k) = n^{-1} \sum_i \left( Y_i - \hat{\beta}_k^T \mathbf{x}_i \right)^2 / (1 - q_{k,ii})^2.$$

If any or some of the values of  $q_{k,ii}$  is/are close to 1, then the cross-validated estimate of  $PE(k)$  may become unstable. For this reason, it has been proposed to replace  $q_{k,ii}$  by its average

$$\bar{q}_k = \sum_i q_{k,ii} / n = \text{trace}(\mathbf{Q}_k) / n = p_k / n.$$

This leads to what is known as the generalized cross-validation estimate

$$\widehat{PE}^{(gcv)}(k) = n^{-1} \sum_i \left( Y_i - \hat{\beta}_k^T \mathbf{x}_i \right)^2 / (1 - p_k / n)^2.$$

It should be noted that the expression of the  $\widehat{PE}^{(gcv)}(k)$  is almost proportional to Akaike's FPE criterion. More precisely,

$$\widehat{PE}^{(gcv)}(k) = n^{-1} [1 + O(p_k / n)^2] FPE_k,$$

assuming that  $p_k / n \rightarrow 0$  as  $n \rightarrow \infty$ .

We now verify the identities given in Eqs. (6a) and (6b). For notational convenience we suppress  $k$  in the expressions for  $\mathbf{X}_k$ ,  $\hat{\beta}_k$ ,  $\mathbf{x}_{k,i}$ , etc. Let  $\mathbf{X}_{(-i)}$  be the matrix obtained from  $\mathbf{X}$  by deleting its  $i$ th row. Thus  $\mathbf{X}_{(-i)}$  has  $n - 1$  rows. Similarly let  $\mathbf{Y}_{(-i)}$  be the vector with

$n - 1$  rows after deleting the  $i$ th row from  $\mathbf{Y}$ . Note that the normal equations for the least squares estimate of  $\boldsymbol{\beta}$  after deleting the  $i$ th observation  $(Y_i, \mathbf{x}_i)$  is given by the solution of  $\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)} \boldsymbol{\beta} = \mathbf{X}_{(-i)}^T \mathbf{Y}_{(-i)}$ . Since

$$\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)} = \mathbf{X}^T \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^T \quad \text{and} \quad \mathbf{X}_{(-i)}^T \mathbf{Y}_{(-i)} = \mathbf{X}^T \mathbf{Y} - \mathbf{x}_i Y_i,$$

using the Sherman-Morrison formula (Section B.1), and denoting  $(\mathbf{X}^T \mathbf{X})^{-1}$  by  $\mathbf{D}$ , we have

$$\begin{aligned} (\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} &= (\mathbf{X}^T \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^T)^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} + \left[ 1 - \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \right]^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} + (1 - q_{ii})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{D} + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i \mathbf{x}_i^T \mathbf{D}, \end{aligned}$$

where the  $i$ th diagonal element of  $\mathbf{Q}$  is  $q_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i = \mathbf{x}_i^T \mathbf{D} \mathbf{x}_i$ . Thus

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{(-i)} &= (\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} \mathbf{X}_{(-i)}^T \mathbf{Y}_{(-i)} \\ &= [\mathbf{D} + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i \mathbf{x}_i^T \mathbf{D}] [\mathbf{X}^T \mathbf{Y} - \mathbf{x}_i Y_i] \\ &= [\mathbf{I} + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i \mathbf{x}_i^T] [\mathbf{D} \mathbf{X}^T \mathbf{Y} - \mathbf{D} \mathbf{x}_i Y_i] \\ &= [\mathbf{I} + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i \mathbf{x}_i^T] [\hat{\boldsymbol{\beta}} - \mathbf{D} \mathbf{x}_i Y_i] \\ &= \hat{\boldsymbol{\beta}} - \mathbf{D} \mathbf{x}_i Y_i + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i [\mathbf{x}_i^T \hat{\boldsymbol{\beta}} - \mathbf{x}_i^T \mathbf{D} \mathbf{x}_i Y_i] \\ &= \hat{\boldsymbol{\beta}} - \mathbf{D} \mathbf{x}_i Y_i + (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i [\hat{Y}_i - q_{ii} Y_i] \\ &= \hat{\boldsymbol{\beta}} - \mathbf{D} \mathbf{x}_i [Y_i - (1 - q_{ii})^{-1} (\hat{Y}_i - q_{ii} Y_i)] \\ &= \hat{\boldsymbol{\beta}} - \mathbf{D} \mathbf{x}_i [(1 - q_{ii})^{-1} (Y_i - \hat{Y}_i)] = \hat{\boldsymbol{\beta}} - (1 - q_{ii})^{-1} \mathbf{D} \mathbf{x}_i \hat{\varepsilon}_i. \end{aligned}$$

This shows that equality (6a) holds. In order to verify Eq. (6b), note that

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{(-i)}^T \mathbf{x}_i &= \hat{\boldsymbol{\beta}}^T \mathbf{x}_i - (1 - q_{ii})^{-1} \mathbf{x}_i^T \mathbf{D} \mathbf{x}_i \hat{\varepsilon}_i \\ &= \hat{\boldsymbol{\beta}}^T \mathbf{x}_i - (1 - q_{ii})^{-1} q_{ii} \hat{\varepsilon}_i, \text{ and} \\ Y_i - \hat{\boldsymbol{\beta}}_{(-i)}^T \mathbf{x}_i &= \left( Y_i - \hat{\boldsymbol{\beta}}^T \mathbf{x}_i \right) + (1 - q_{ii})^{-1} q_{ii} \hat{\varepsilon}_i \\ &= \hat{\varepsilon}_i + (1 - q_{ii})^{-1} q_{ii} \hat{\varepsilon}_i = (1 - q_{ii})^{-1} \hat{\varepsilon}_i. \end{aligned}$$

## 11.9 Some Alternate Methods for Regression

In a regression setting with  $p - 1$  independent variables, in some cases,  $Y$  is related to only  $k < p - 1$  independent variables, and thus it makes sense to remove some of the independent variables from the model. Stepwise regression, all subsets regression, or a penalty method such as lasso may be used for this purpose. In some cases, another situation may be true where  $Y$  is related to all the independent variables and the independent variables are well correlated among themselves. In such cases, it may be reasonable to regress  $Y$  on a few (say  $k < p - 1$ ) appropriately created linear combinations of the independent variables. Partial least squares (PLS) or principal components regression (PCR) may be used for this purpose. The phenomenon of all or some of the independent variables being well-correlated among themselves is called *multicollinearity*. Another method that has been proposed in the literature for multicollinear cases is known as the ridge regression.

Theoretical properties of some of the procedures described above are either nontrivial or not well-understood. For this reason, only outlines of these procedures will be given. It is also important to point out there are other regression procedures in addition to the ones mentioned above.

For all the methods to be discussed here, we assume the Gauss-Markov setup given in Eq. (1) is true with  $\epsilon$  as the vector of iid variables with mean 0 and variance  $\sigma^2$ , and the goal is to estimate the mean vector  $X\beta$  on the basis of the observation vector  $Y$ . Even though the typical assumption is that  $\{\varepsilon_i\}$  are iid, in some cases, such as stepwise regression, the typical assumption is  $\epsilon \sim N_n(\mathbf{0}, \sigma^2 I)$ .

### 11.9.1 All Subsets Regression

All subsets regression consists of fitting all possible submodels and choosing an appropriate (sub)model using a criterion such as AIC, Mallows' or BIC (ie, choose the model with the smallest value of the criterion). If there are  $p - 1$  independent variables then there are  $2^{p-1}$  possible submodels. When  $p$  is not small, then the number of submodels  $2^{p-1}$  is quite large and, in such cases, carrying out all subsets regression may not be feasible. For this reason, one may use a computationally feasible method such as stepwise regression which is described below.

### 11.9.2 Stepwise Regression

Let us first discuss the simplest versions of this method: forward selection and backward elimination.

In the forward selection method, one builds up a model beginning with no independent variable and adding one variable at a time till it is not possible to add any more. It is typical to start with the model  $Y = \beta_0 + \varepsilon$  and then try each of the  $p - 1$  regression models,  $Y = \beta_0 + \beta_j X_j + \varepsilon, j = 1, \dots, p - 1$ . When comparing the models  $Y = \beta_0 + \varepsilon$  (reduced)

to  $Y = \beta_0 + \beta_j X_j + \varepsilon$  (full), one can test  $H_0: \beta_j = 0$  against  $H_1: \beta_j \neq 0$ . This leads to an  $F$ -statistic and, since there are  $p - 1$  models, there are  $p - 1$  such  $F$ -values. The variable with the largest  $F$ -value is the best candidate for inclusion. If this largest  $F$ -statistic is larger than a preselected threshold (called  $F$ -to-enter), then the corresponding independent variable is entered in the model. If the value of the largest  $F$ -statistic is smaller than the threshold, then no variable can be added and it is declared that  $Y = \beta_0 + \varepsilon$  is the most appropriate model.

Suppose that variable  $X_2$  has been selected in step 1, then one considers  $p - 2$  models,  $Y = \beta_0 + \beta_2 X_2 + \beta_j X_j + \varepsilon$ ,  $1 \leq j \neq 2 \leq p - 1$ . Now there are  $p - 2$  testing problems  $H_0: \beta_j = 0$  against  $H_1: \beta_j \neq 0$ ,  $1 \leq j \neq 2 \leq p - 1$ . Thus there are  $p - 2$   $F$ -values and the variable with the largest  $F$ -value is entered if it is larger than the threshold. Otherwise  $Y = \beta_0 + \beta_2 X_2 + \varepsilon$  is considered the most appropriate model. In this manner, independent variables can be added one at a time till it is not possible to enter any more.

*Remark 11.9.1.* In each step, instead of using the  $F$ -values, one may use  $p$ -values and enter a variable if the corresponding  $p$ -value is the smallest, and it is smaller than a threshold  $p$ -value (called  $p$ -to-enter).

The selected model depends on the threshold  $F$ -value (often taken to be equal to 3.5 or 4) or the threshold  $p$ -value (often taken to be equal to 0.05 or 0.10). Changing the threshold value may lead to a different model. For this reason, sometimes it is considered appropriate to look at the best candidate's AIC or BIC value at each step. The model with the smallest value AIC or BIC is considered the most appropriate.

In the backward elimination method, one starts with the full model and then deletes variables one at a time using the  $F$ -values and a threshold called  $F$ -to-delete till no variable can be dropped. As in the forward selection method, one may use  $p$ -values instead of  $F$ -values or use a criterion such as AIC to carry out this procedure.

Stepwise regression can be done in a forward or in a backward manner. In the forward stepwise regression, variables are added as in the forward stepwise regression. However, in each step it also has the option of deleting a variable that is already in the model. So forward stepwise regression is basically a forward selection method which includes a backward elimination step. For this reason, a forward stepwise method needs two threshold values: one for forward selection and the other for backward elimination.

Similarly, backward stepwise regression is basically a backward elimination method with the option of reentering a variable that is outside the model. As in the forward stepwise method, it needs two threshold values.

### 11.9.3 Penalty Methods

Penalty-based approaches to regression seek to minimize the least squares criterion  $G(\mathbf{b}) = \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|^2$  with restrictions on  $\mathbf{b}$ . This can usually be achieved by minimizing  $G(\mathbf{b})$  plus a penalty term as given below

$$G(\mathbf{b}, \lambda) = G(\mathbf{b}) + \lambda p(\mathbf{b}),$$

where the penalty parameter  $\lambda$  is nonnegative and the penalty function  $p(\mathbf{b}) \geq 0$  is equal to 0 at  $\mathbf{b} = \mathbf{0}$ . There are many possible choices for the penalty function, but we will mention only two:  $p(\mathbf{b}) = \sum c_j |b_j|$  and  $p(\mathbf{b}) = \sum c_j^2 b_j^2$ , where  $c_j^2$  is the  $j$ th diagonal element of  $\mathbf{X}^T \mathbf{X}$ . For the first case, also known as lasso, there is no explicit expression for the minimizer. However, for the second case, known as the ridge regression, there is an explicit expression for the minimizer  $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{C})^{-1} \mathbf{X}^T \mathbf{Y}$ , where  $\mathbf{C}$  is a diagonal matrix with diagonal elements  $\{c_j^2\}$ . When the penalty term  $\lambda = 0$ , both methods are the same as the ordinary least squares. In general, the penalty parameter can be chosen by cross-validation or its modified versions.

Ridge regression works well in estimating the mean response when the independent variables are well correlated among themselves. However, if the independent variables are mutually uncorrelated, its performance may not be satisfactory. With appropriate choice of  $\lambda$  and under regularity conditions, lasso can recover the regression model when some of the beta parameters are nonzero and the rest are zeros. However, the ridge estimates will typically result in nonzero estimates for all the beta parameters.

The penalty methods mentioned above can be motivated by a Bayesian consideration where the true parameters are assumed to randomly distributed. If all the variables have been standardized, then there is no need for an intercept term and the penalty term can be taken to be equal to  $\lambda \sum \beta_j^2 = \lambda \|\boldsymbol{\beta}\|^2$ . If the true regression coefficients  $\{\beta_j\}$  are assumed to be iid  $N(0, \tau^2)$ , then ignoring the terms that do not depend on  $\mathbf{Y}$  or  $\boldsymbol{\beta}$ ,  $-2$  times the logarithm of the joint pdf of  $\mathbf{Y}$  and  $\boldsymbol{\beta}$  is

$$\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2/\sigma^2 + \|\boldsymbol{\beta}\|^2/\tau^2 = \sigma^{-2} \left[ \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \right], \quad \lambda = \sigma^2/\tau^2.$$

Minimizing the above with respect to  $\boldsymbol{\beta}$  leads to the ridge estimate of  $\boldsymbol{\beta}$ . If instead of assuming normal distribution of  $\{\beta_j\}$ , one assumes  $\{\beta_j\}$  to be iid double exponential with the pdf  $(2\tau)^{-1} \exp[-|u|/\tau]$ ,  $-\infty < u < \infty$ ,  $\tau > 0$ , then the argument above leads to a lasso estimate of  $\boldsymbol{\beta}$ .

It can be shown that there exists  $\lambda > 0$  such that the ridge estimate  $\hat{\mathbf{X}}\hat{\boldsymbol{\beta}}(\lambda)$  is a superior estimate of  $\mu = E[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}$  than the least squares estimate  $\hat{\mathbf{X}}\hat{\boldsymbol{\beta}}$  of  $\mu$ . We write down the result for the case when all the variables have been standardized so that the penalty term is  $\lambda \|\boldsymbol{\beta}\|^2$  and the ridge estimate of  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ .

**Lemma 11.9.1.** *Let us assume that all the variables have been standardized in a regression model and the ridge estimate of  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ ,  $\lambda > 0$ . Let  $D(\lambda) = E[\|\mathbf{X}\hat{\boldsymbol{\beta}}(\lambda) - \mathbf{X}\boldsymbol{\beta}\|^2]$  be the expected value of the squared distance between  $\mathbf{X}\hat{\boldsymbol{\beta}}(\lambda)$  and  $\mathbf{X}\boldsymbol{\beta}$ . There exists  $\lambda > 0$  such that  $D(\lambda) < D(0)$ , ie, there exists a ridge estimate  $\mathbf{X}\hat{\boldsymbol{\beta}}(\lambda)$ ,  $\lambda > 0$ , which is a superior estimate of  $\mathbf{X}\boldsymbol{\beta}$  than the least squares estimate  $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\hat{\boldsymbol{\beta}}(0)$ .*

*Proof of Lemma 11.9.1.* Clearly,

$$\boldsymbol{\beta}(\lambda) = E[\hat{\boldsymbol{\beta}}(\lambda)] = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta} - \lambda (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\beta}.$$

Hence, writing  $\mathbf{A}_\lambda = \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  for notational convenience, we have

$$\begin{aligned}\mathbf{X}\hat{\boldsymbol{\beta}}(\lambda) - \mathbf{X}\boldsymbol{\beta} &= \mathbf{X}\hat{\boldsymbol{\beta}}(\lambda) - \mathbf{X}\boldsymbol{\beta}(\lambda) + \mathbf{X}\boldsymbol{\beta}(\lambda) - \mathbf{X}\boldsymbol{\beta} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} - \lambda \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\beta} \\ &= \mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} - \lambda \mathbf{X}\mathbf{A}_\lambda^{-1} \boldsymbol{\beta}, \text{ and} \\ D(\lambda) &= E[\|\mathbf{X}\hat{\boldsymbol{\beta}}(\lambda) - \mathbf{X}\boldsymbol{\beta}\|^2] = E[\|\mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}\|^2] + \lambda^2 \|\mathbf{X}\mathbf{A}_\lambda^{-1} \boldsymbol{\beta}\|^2 \\ &:= D_V(\lambda) + D_B(\lambda).\end{aligned}$$

We may regard  $D_V(\lambda)$  and  $D_B(\lambda)$  as the variance and the bias-squared terms. Since  $\text{Cov}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}$ , we have

$$\begin{aligned}D_V(\lambda) &= E[\|\mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}\|^2] = \text{trace}(\mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \text{Cov}[\boldsymbol{\varepsilon}] \mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T) \\ &= \sigma^2 \text{trace}(\mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T) \\ &= \sigma^2 \text{trace}(\mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1} \mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1}) = \sigma^2 \text{trace}([\mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1}]^2).\end{aligned}$$

Let the spectral decomposition of  $\mathbf{X}^T \mathbf{X}$  be  $\sum_{j=1}^p \lambda_j \mathbf{u}_j \mathbf{u}_j^T$  with  $\lambda_1 \geq \dots \geq \lambda_p$ , and denote  $\gamma_j = \mathbf{u}_j^T \boldsymbol{\beta}$ . Then

$$\begin{aligned}\mathbf{A}_\lambda^{-1} &= \sum (\lambda_j + \lambda)^{-1} \mathbf{u}_j \mathbf{u}_j^T, \\ \mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1} &= \sum \lambda_j (\lambda_j + \lambda)^{-1} \mathbf{u}_j \mathbf{u}_j^T, \text{ and} \\ \mathbf{A}_\lambda^{-1} \mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1} &= \sum \lambda_j (\lambda_j + \lambda)^{-2} \mathbf{u}_j \mathbf{u}_j^T.\end{aligned}$$

Therefore

$$\begin{aligned}D_V(\lambda) &= \sigma^2 \text{trace}([\mathbf{X}^T \mathbf{X}\mathbf{A}_\lambda^{-1}]^2) = \sum [\lambda_j / (\lambda_j + \lambda)]^2, \\ D_B(\lambda) &= \lambda^2 \|\mathbf{X}\mathbf{A}_\lambda^{-1} \boldsymbol{\beta}\|^2 = \lambda^2 \sum \lambda_j (\lambda_j + \lambda)^{-2} \gamma_j^2.\end{aligned}$$

It is fairly easy to see that  $D_V(\lambda)$  is a decreasing function of  $\lambda$ , but  $D_B(\lambda)$  is increasing in  $\lambda$ . As a matter of fact

$$\begin{aligned}D'_V(\lambda) &= -2\sigma^2 \sum \lambda_j^2 (\lambda_j + \lambda)^{-3}, \quad D'_B(\lambda) = 2\lambda \sum \lambda_j^{-2} (\lambda_j + \lambda)^{-3} \gamma_j^2, \\ D'(\lambda) &= D'_V(\lambda) + D'_B(\lambda) \\ &= -2\sigma^2 \sum \lambda_j^2 (\lambda_j + \lambda)^{-3} + 2\lambda \sum \lambda_j^{-2} (\lambda_j + \lambda)^{-3} \gamma_j^2, \text{ and} \\ D'(0) &= -2\sigma^2 \sum \lambda_j^{-1} < 0.\end{aligned}$$

Thus  $D(\lambda)$  is decreasing in a neighborhood of 0 and hence  $D(\lambda) < D(0)$  for some  $\lambda > 0$ .  $\square$

### 11.9.4 Partial Least Squares and Principal Components Regression

The main idea behind PCR and PLS is to create a design matrix  $\mathbf{Z}$  from  $\mathbf{X}$  whose columns are mutually orthogonal. Thus if an  $n \times p$  matrix  $\mathbf{Z} = [\mathbf{X}\mathbf{u}_1, \dots, \mathbf{X}\mathbf{u}_p]$  is created so that

$$\mathbf{Z}_j^T \mathbf{Z}_k = (\mathbf{X}\mathbf{u}_j)^T (\mathbf{X}\mathbf{u}_k) = \mathbf{u}_j^T \mathbf{X}^T \mathbf{X}\mathbf{u}_k = 0, \quad j \neq k,$$

then  $\mathbf{Z}^T \mathbf{Z}$  is a diagonal matrix. If we write  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$ , then  $\mathbf{X}\beta = \mathbf{X}\mathbf{U}\mathbf{U}^{-1}\beta = \mathbf{Z}\alpha$ , with  $\alpha = \mathbf{U}^{-1}\beta$ . So we have a new Gauss-Markov model  $\mathbf{Y} = \mathbf{Z}\alpha + \boldsymbol{\epsilon}$  and the least squares estimate of  $\alpha$  is  $\hat{\alpha} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Y}$ , and  $\hat{\alpha}_j = \mathbf{Z}_j^T \mathbf{Y} / \mathbf{Z}_j^T \mathbf{Z}_j$  is the estimate of  $\alpha_j$ , the  $j$ th component of  $\alpha$ . It is important to note that in this formulation, estimate of  $\alpha_j$  depends only on  $\mathbf{Z}_j$  and  $\mathbf{Y}$ ,  $\text{Cov}[\hat{\alpha}] = \sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1}$ , and whenever  $j \neq k$ ,  $\text{Cov}[\hat{\alpha}_j, \hat{\alpha}_k] = 0$  since  $\mathbf{Z}^T \mathbf{Z}$  is a diagonal matrix. Model selection now involves choosing a few of the newly created independent variables  $\mathbf{Z}_1, \dots, \mathbf{Z}_p$  and this can be carried out by using stepwise regression with an appropriate model selection criterion described in the last section. Since  $\beta = \mathbf{U}\alpha$ , once we have an estimate  $\hat{\alpha}$  of  $\alpha$  (and this includes the case when some  $\alpha_j$  are set be zero if we decide to delete the corresponding  $\mathbf{Z}_j$ 's), we can get the estimate of  $\beta$  as  $\hat{\beta} = \mathbf{U}\hat{\alpha} = \sum \hat{\alpha}_j \mathbf{u}_j$ .

PCR and PLS differ in the way the vectors  $\mathbf{u}_1, \dots, \mathbf{u}_p$  are created or the new mutually orthogonal vectors  $\mathbf{Z}_1, \dots, \mathbf{Z}_p$  are created. In PCR,  $\{\mathbf{u}_j\}$  are taken to be the orthonormal eigenvectors obtained from the spectral decomposition of  $\mathbf{X}^T \mathbf{X}$ . Thus if  $\lambda_1, \dots, \lambda_p$  are the eigenvalues of  $\mathbf{X}^T \mathbf{X}$  with the corresponding orthonormal eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_p$ , then  $\mathbf{Z}_j$  is taken to be  $\mathbf{Z}_j = \mathbf{X}\mathbf{u}_j, j = 1, \dots, p$ , then clearly  $\{\mathbf{Z}_j\}$  are mutually orthogonal.

In the presence of multicollinearity, it is of interest to create the mutually uncorrelated variables economically so that a linear combination of a few of them can model the response variable  $\mathbf{Y}$  well. It is thus desirable to have a method which creates the vectors  $\{\mathbf{u}_j\}$  using both  $\mathbf{Y}$  and  $\mathbf{X}$ . PCR uses only information from  $\mathbf{X}$ , whereas PLS uses both  $\mathbf{X}$  and  $\mathbf{Y}$ . It should be pointed out that PLS is basically what is known as the Conjugate Gradient Method in Numerical Analysis. It is an iterative method for obtaining the minimum of a quadratic form  $\beta^T \mathbf{S} \beta - 2\mathbf{b}^T \beta$  over  $\beta \in \mathbb{R}^p$ , where  $\mathbf{S}$  is a  $p \times p$  positive definite matrix.

*Remark 11.9.2.* It is generally advisable to standardize all the variables for the PCR and the PLS so that there is no need to include an intercept term in the model and all the variables are on the same scale.

There are many equivalent ways of describing PLS, including a version with statistical interpretation. This procedure generates  $\mathbf{u}_1, \mathbf{u}_2, \dots$  in a recursive manner along with the estimated coefficients  $\hat{\alpha}_1, \hat{\alpha}_2, \dots$  and the corresponding estimates of  $\beta$ . Here is a version which follows the numerical analyst's description. For notational convenience, let  $\mathbf{S} = n^{-1} \mathbf{X}^T \mathbf{X}$ ,  $\mathbf{b} = n^{-1} \mathbf{X}^T \mathbf{Y}$  so that the normal equations can be written as  $\mathbf{S}\beta = \mathbf{b}$ . The following iterative scheme describes how the columns of the matrix  $\mathbf{Z}$ , and  $\hat{\alpha}_1, \hat{\alpha}_2, \dots$  etc., are generated recursively.

*Step 1.* Take  $\mathbf{u}_1 = \mathbf{b}$ . Then  $\mathbf{Z}_1 = \mathbf{X}\mathbf{u}_1$ ,  $\hat{\alpha}_1 = \mathbf{Z}_1^T \mathbf{Y} / \mathbf{Z}_1^T \mathbf{Z}_1 = \mathbf{u}_1^T \mathbf{b} / \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ , and  $\hat{\beta}^{(1)} = \hat{\alpha}_1 \mathbf{u}_1$ .

Denote the remainder  $\mathbf{b} - \mathbf{S}\hat{\beta}^{(1)}$  by  $\mathbf{r}_1$ .

We will describe how to carry out this recursion. Suppose that  $\mathbf{u}_1, \dots, \mathbf{u}_k$  has been generated with the corresponding estimates  $\hat{\beta}^{(1)}, \dots, \hat{\beta}^{(k)}$  and the remainders

$\mathbf{r}_1 = \mathbf{b} - \mathbf{S}\hat{\beta}^{(1)}, \dots, \mathbf{r}_k = \mathbf{b} - \mathbf{S}\hat{\beta}^{(k)}$ , and we now want to create  $\mathbf{u}_{k+1}$ .

Step  $k+1$ . Let  $\mathbf{u}_{k+1} = \mathbf{r}_k - c_k \mathbf{u}_k$ , where  $c_k$  is such that  $\mathbf{Z}_{k+1} = \mathbf{X}\mathbf{u}_{k+1}$  is orthogonal to  $\mathbf{Z}_k = \mathbf{X}\mathbf{u}_k$ . This leads to  $c_k = \mathbf{u}_k^T \mathbf{S} \mathbf{r}_k / \mathbf{u}_k^T \mathbf{S} \mathbf{u}_k$ . So we have

$$\mathbf{u}_{k+1} = \mathbf{r}_k - c_k \mathbf{u}_k, \quad \mathbf{Z}_{k+1} = \mathbf{X}\mathbf{u}_{k+1} \text{ with } c_k = \mathbf{u}_k^T \mathbf{S} \mathbf{r}_k / \mathbf{u}_k^T \mathbf{S} \mathbf{u}_k,$$

$$\hat{\alpha}_{k+1} = \mathbf{Z}_{k+1}^T \mathbf{Y} / \mathbf{Z}_{k+1}^T \mathbf{Z}_{k+1} = \mathbf{u}_{k+1}^T \mathbf{b} / \mathbf{u}_{k+1}^T \mathbf{S} \mathbf{u}_{k+1}, \text{ and}$$

$$\hat{\beta}^{(k+1)} = \hat{\alpha}_1 \mathbf{u}_1 + \dots + \hat{\alpha}_k \mathbf{u}_k + \hat{\alpha}_{k+1} \mathbf{u}_{k+1} = \hat{\beta}^{(k)} + \hat{\alpha}_{k+1} \mathbf{u}_{k+1},$$

and the remainder is  $\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{X}\hat{\beta}^{(k+1)}$ .

*Remark 11.9.3.*

- (a) It can be shown that for the procedure described above  $\mathbf{Z}_1 = \mathbf{X}\mathbf{u}_1, \mathbf{Z}_2 = \mathbf{X}\mathbf{u}_2, \dots$  are orthogonal.
- (b) If  $\mathbf{S}$  is a diagonal matrix, the iteration stops after the first iteration. If  $\mathbf{X}$  is  $n \times p$ , then the procedure stops after  $p$  iterations and  $\hat{\beta}^{(p)}$  equals the least squares estimate  $\hat{\beta}$  of  $\beta$ .
- (c) How do we decide when to stop the iterations? One may use a criterion such as AIC or cross-validation and terminate the iteration when the value of the AIC stops decreasing.

## 11.10 Random- and Mixed-Effects Models

Random-effects models come up in many situations of practical interest. Consider a simple example. Suppose in a state, it is desired to know the average performance of children in some standardized mathematics test. Since the state has many schools and the student performance may vary from school to school, it may be desirable to choose  $k$  schools at random and, for each selected school, the standardized test is given to some randomly selected children. Thus if  $Y_{ij}$  is the performance of the  $j$ th child in the  $i$ th school, one may write a one-factor ANOVA model as  $Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$ , where  $\mu$  is the overall mean,  $\alpha_i$  is the school effect, and  $\varepsilon_{ij}$  is the random error. Since the schools are randomly selected it is reasonable to assume that, for this one-factor model,  $\{\alpha_i\}$  are iid  $N(0, \sigma_1^2)$ . This is perhaps the simplest of all random-effect models.

**Example 11.10.1** (One random factor). Suppose we have observations  $\{Y_{ij} : j = 1, \dots, n_i, i = 1, \dots, k\}$  where the random factor has  $k$  levels and the model is

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\mu$  is a constant,

- (i)  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ ,  $\{\alpha_i\}$  are iid  $N(0, \sigma_1^2)$ , and  
(ii)  $\{\alpha_i\}$  and  $\{\varepsilon_{ij}\}$  are mutually independent.

Thus, we have

$$\begin{aligned} E[Y_{ij}] &= \mu, \quad \text{Var}[Y_{ij}] = \sigma_1^2 + \sigma^2, \\ \text{Cov}[Y_{ij}, Y_{ij'}] &= \sigma_1^2, \quad j \neq j', \text{ and} \\ \text{Cov}[Y_{ij}, Y_{i'j'}] &= 0, \quad i \neq i'. \end{aligned}$$

Now the variability of  $Y_{ij}$  has two components: variability among schools and the error variance. For this reason,  $\sigma_1^2$  and  $\sigma^2$  are also called the components of variance of  $Y_{ij}$ . The observations  $\{Y_{i1}, \dots, Y_{in_i}\}$  in a particular school are correlated, unlike in the fixed factor case. The correlation between  $Y_{ij}$  and  $Y_{ij'}, j \neq j'$ , which equals  $\sigma_1^2 / (\sigma_1^2 + \sigma^2)$ , is called the intraclass correlation. The goal is often to estimate  $\mu$  and the intraclass correlation coefficient.

**Example 11.10.2** (Two Random Factors). Suppose we have observations  $\{Y_{ijk}: k = 1, \dots, n_{ij}, j = 1, \dots, b, i = 1, \dots, a\}$  where factor  $A$  has  $a$  levels, factor  $B$  has  $b$  levels, and both factors are random. Thus it can be written as

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk},$$

where it is assumed that  $\mu$  is a constant;  $\{\alpha_i\}$  are iid  $N(0, \sigma_1^2)$ ;  $\{\beta_j\}$  are iid  $N(0, \sigma_2^2)$ ;  $\{(\alpha\beta)_{ij}\}$  are iid  $N(0, \sigma_3^2)$ ;  $\{\varepsilon_{ijk}\}$  are iid  $N(0, \sigma^2)$ ; and  $\{\alpha_i\}$ ,  $\{\beta_j\}$ ,  $\{(\alpha\beta)_{ij}\}$ , and  $\{\varepsilon_{ijk}\}$ , are mutually independent. Here

$$E[Y_{ijk}] = \mu \text{ and } \text{Var}[Y_{ijk}] = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \sigma^2.$$

So the components of variance are  $\sigma_1^2$ ,  $\sigma_2^2$ ,  $\sigma_3^2$ , and  $\sigma^2$ . The goal is often to estimate these variance components along with the overall mean  $\mu$ .

**Example 11.10.3** (Mixed-Effects Model: One Factor Fixed and One Factor Random). In the last example, if factor  $A$  is fixed but factor  $B$  is random, then the model is the same except that  $\mu$  and  $\alpha_i$  are constants with  $\sum \alpha_i = 0$ , but  $\{\beta_j\}$ ,  $\{(\alpha\beta)_{ij}\}$  are  $\{\varepsilon_{ijk}\}$  are iid  $N(0, \sigma_1^2)$ ,  $N(0, \sigma_2^2)$  and  $N(0, \sigma^2)$ , respectively. It is also assumed that  $\{\beta_j\}$ ,  $\{(\alpha\beta)_{ij}\}$  and  $\{\varepsilon_{ijk}\}$  are mutually independent. Thus

$$E[Y_{ij}] = \mu + \alpha_i \quad \text{and} \quad \text{Var}[Y_{ij}] = \sigma_1^2 + \sigma_2^2 + \sigma^2.$$

This is a mixed-effects model with the variance components  $\sigma_1^2$ ,  $\sigma_2^2$ , and  $\sigma^2$ .

**Example 11.10.4** (Mixed-Effects Model: One Random Factor With a Covariate). We can also have a mixed model in which some factors are random, some fixed, and there are one or more covariates. Consider a case with one random factor and one covariate. Thus

if  $(Y_{ij}, X_{ij}), j = 1, \dots, n_i$ , are the observed values of the response and the covariate for the factor at level  $i$ , then a model can be written as

$$Y_{ij} = \beta_0 + \beta_1 X_{ij} + \alpha_i + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\beta_0, \beta_1$  are unknown constants,  $\{\alpha_i\}$  are iid  $N(0, \sigma_1^2)$ ,  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ , and  $\{\alpha_i\}$  and  $\{\varepsilon_{ij}\}$  are mutually independent. In this case, we have

$$\mathbb{E}[Y_{ij}] = \beta_0 + \beta_1 X_{ij} \quad \text{and} \quad \text{Var}[Y_{ij}] = \sigma_1^2 + \sigma^2.$$

**Example 11.10.5** (Mixed Model: Growth Model With Random Slopes). Suppose we have  $k$  children whose heights (growths) are measured at various ages. Let  $Y_{ij}$  be the height of the  $i$ th child at age  $t_j$ ,  $j = 1, \dots, n_0$ , and let  $X_{ij}$  be the vector of covariates (available nutrition, parents' heights, family income, etc.). If we model the height as a polynomial of degree  $r$  in age with random slopes, then a reasonable model is

$$Y_{ij} = \beta_0 + \beta_1^T X_{ij} + \sum_{l=1}^r \gamma_{il} t_j^l + \varepsilon_{ij}, \quad j = 1, \dots, n_0, \quad i = 1, \dots, k,$$

where  $\beta_0, \beta_1$  are nonrandom,  $\{\gamma_{il}: i = 1, \dots, k\}$  are iid  $N(0, \sigma_l^2)$  for each  $l$ ,  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ , and  $\{\gamma_{il}\}$  and  $\{\varepsilon_{ij}\}$  are mutually independent. This is an example of a growth model which allows different rates of growths for different children.

In Example 11.10.4 we can create  $k$  indicator variables  $\{Z_{ij1}, \dots, Z_{ijk}\}$  as follows. Let  $Z_{ijl} = 1$  if  $i = l$  and 0 otherwise. If we denote the random effects by  $\{\gamma_i\}$  instead of  $\{\alpha_i\}$ , then

$$Y_{ij} = \beta_0 + \beta_1 X_{ij} + Z_{ij1} \gamma_1 + \dots + Z_{ijk} \gamma_k + \varepsilon_{ij}.$$

We can express this model in a matrix form. Let  $n = n_1 + \dots + n_k$  be the total number of observations. Let  $\mathbf{X}$  be the  $n \times 2$  matrix whose first column has only ones and its second column consists of values of the covariate. Let  $\mathbf{Z}$  be an  $n \times k$  matrix whose first column consists of the values of  $Z_{ij1}$ , second column consists of the values  $Z_{ij2}$ , etc. Then we may rewrite the last model as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon},$$

where  $\boldsymbol{\gamma}$  is the  $k \times 1$  vector of  $\gamma_1, \dots, \gamma_k$ . Here  $\mathbf{X}$  and  $\mathbf{Z}$  are known matrices,  $\boldsymbol{\beta}$  is the vector of unknown parameters,  $\boldsymbol{\gamma} \sim N_k(\mathbf{0}, \sigma_1^2 \mathbf{I})$ ,  $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$ , and  $\boldsymbol{\gamma}$  and  $\boldsymbol{\varepsilon}$  are independent.

One-factor random-effects model can also be written in this form. In order to accommodate the two-factor random- and mixed-effects models, which have more than two variance components, we consider a more general model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1 \boldsymbol{\gamma}_1 + \dots + \mathbf{Z}_r \boldsymbol{\gamma}_r + \boldsymbol{\varepsilon}, \tag{7}$$

where  $\mathbf{X}, \mathbf{Z}_1, \dots, \mathbf{Z}_r$  are known matrices of orders  $n \times p, n \times q_1, \dots, n \times q_r$ , respectively,  $\boldsymbol{\beta}$  is a vector of unknown parameters,  $\boldsymbol{\gamma}_j \sim N_{q_j}(\mathbf{0}, \sigma_j^2 \mathbf{I}), j = 1, \dots, r$ ,  $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$ , and  $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_k$  and  $\boldsymbol{\varepsilon}$  are all independent. In this framework

$$\mathbf{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta} \text{ and } \text{Cov}[\mathbf{Y}] = \sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1^T + \dots + \sigma_k^2 \mathbf{Z}_r \mathbf{Z}_r^T + \sigma^2 \mathbf{I}.$$

Two-factor mixed- and random-effects models (Examples 11.10.2 and 11.10.5) can also be expressed in this form and it is left as an exercise. It is often useful to express the above model by taking  $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$  and  $\boldsymbol{\gamma}$  as a column vector with  $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_k$  stacked vertically. Thus  $\mathbf{Z}$  is  $n \times q$  and  $\boldsymbol{\gamma}$  is  $q \times 1$ , where  $q = q_1 + \dots + q_k$ . Thus

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}, \quad (8)$$

where  $\boldsymbol{\gamma} \sim N_m(\mathbf{0}, \mathbf{D})$ , where  $\mathbf{D}$  is a block diagonal matrix with the diagonal blocks  $\text{Cov}[\boldsymbol{\gamma}_1], \dots, \text{Cov}[\boldsymbol{\gamma}_r]$ .

## Assumption

For the model in Eq. (8) we assume that  $\text{rank}(\mathbf{X}) = p, \text{rank}(\mathbf{Z}_j) = q_j, j = 1, \dots, r$ .

### 11.10.1 Estimation of $\boldsymbol{\beta}$ and Prediction of $\boldsymbol{\gamma}$

We now discuss estimation of  $\boldsymbol{\beta}$  and prediction of  $\boldsymbol{\gamma}$  assuming all the variance components to be known. Since  $\boldsymbol{\gamma}$  is random, its estimation (or linear function of it) is called prediction.

#### Definition 11.10.1.

- (i) A linear function  $\hat{\boldsymbol{\beta}}$  of  $\mathbf{Y}$  is called a BLUE of  $\boldsymbol{\beta}$  if it is an unbiased estimator of  $\boldsymbol{\beta}$ , that is,  $\mathbf{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ , and, for any  $\mathbf{a} \in \mathbb{R}^p$ ,  $\text{Var}[\mathbf{a}^T \hat{\boldsymbol{\beta}}] \leq \text{Var}[\mathbf{l}^T \mathbf{Y}]$  for all linear unbiased estimators  $\mathbf{l}^T \mathbf{Y}$  of  $\mathbf{a}^T \boldsymbol{\beta}$ .
- (ii) A linear function  $\hat{\boldsymbol{\gamma}}$  of  $\mathbf{Y}$  is called a best linear unbiased predictor (BLUP) of  $\boldsymbol{\gamma}$  if it is an unbiased predictor of  $\boldsymbol{\gamma}$ , that is,  $\mathbf{E}[\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}] = \mathbf{0}$ , and for any  $\mathbf{b} \in \mathbb{R}^q$ ,  $\mathbf{E}[\mathbf{b}^T \hat{\boldsymbol{\gamma}} - \mathbf{b}^T \boldsymbol{\gamma}]^2 \leq \mathbf{E}[\mathbf{l}^T \mathbf{Y} - \mathbf{b}^T \boldsymbol{\gamma}]^2$  for all linear unbiased predictors  $\mathbf{l}^T \mathbf{Y}$  of  $\mathbf{b}^T \boldsymbol{\gamma}$ .

According to the definition given above,  $\hat{\boldsymbol{\beta}}$  is a BLUE of  $\boldsymbol{\beta}$  if  $\mathbf{a}^T \hat{\boldsymbol{\beta}}$  is the BLUE of  $\mathbf{a}^T \boldsymbol{\beta}$  for any  $\mathbf{a} \in \mathbb{R}^p$ . Similarly,  $\hat{\boldsymbol{\gamma}}$  is BLUP for  $\boldsymbol{\gamma}$  if  $\mathbf{b}^T \hat{\boldsymbol{\gamma}}$  is the BLUP of  $\mathbf{b}^T \boldsymbol{\gamma}$  for any  $\mathbf{b} \in \mathbb{R}^q$ .

For the mixed model, the BLUE of  $\boldsymbol{\beta}$  is no longer equal to  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . As a matter of fact, it now depends on the unknown matrix  $\mathbf{D}$  and  $\sigma^2$ . Let  $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{Y}] = \mathbf{Z} \mathbf{D} \mathbf{Z}^T + \sigma^2 \mathbf{I}$ . Then, if we multiply both sides of the mixed model by  $\boldsymbol{\Sigma}^{-1/2}$ , where  $\boldsymbol{\Sigma}^{1/2}$  is a symmetric square root of  $\boldsymbol{\Sigma}$ , it leads to a reexpression of the mixed model as  $\tilde{\mathbf{Y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}}$ , where  $\tilde{\mathbf{Y}} = \boldsymbol{\Sigma}^{-1/2} \mathbf{Y}, \tilde{\mathbf{X}} = \boldsymbol{\Sigma}^{-1/2} \mathbf{X}$ , and  $\tilde{\boldsymbol{\varepsilon}} = \boldsymbol{\Sigma}^{-1/2} (\mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon})$ . Thus the reexpressed mixed model is in the standard Gauss-Markov setup since  $\mathbf{E}[\tilde{\boldsymbol{\varepsilon}}] = \mathbf{0}$  and  $\text{Cov}[\tilde{\boldsymbol{\varepsilon}}] = \mathbf{I}$ , and the BLUE of  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{Y}} = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y}$ . We leave it to the reader to verify that the BLUE of  $\boldsymbol{\beta}$  is unique.

Let us first find out what the BLUP of  $\mathbf{b}^T \boldsymbol{\gamma}$  would be if  $\boldsymbol{\beta}$  were known. The best linear predictor of  $\mathbf{b}^T \boldsymbol{\gamma}$  can be obtained by a linear regression of  $\mathbf{b}^T \boldsymbol{\gamma}$  on  $\mathbf{Y}$ . So the best linear predictor  $h_0 + \mathbf{h}_1^T \mathbf{Y}$  has the form

$$\begin{aligned}\mathbf{h}_1 &= [\text{Cov}(\mathbf{Y})]^{-1} \text{Cov}[\mathbf{Y}, \mathbf{b}^T \boldsymbol{\gamma}] = \boldsymbol{\Sigma}^{-1} \mathbf{Z} \mathbf{D} \mathbf{b}, \\ h_0 &= \mathbb{E}[\mathbf{b}^T \boldsymbol{\gamma}] - \mathbf{h}_1^T \mathbb{E}[\mathbf{Y}] = -\mathbf{h}_1^T \mathbf{X} \boldsymbol{\beta}, \text{ and thus} \\ h_0 + \mathbf{h}_1^T \mathbf{Y} &= \mathbf{h}_1^T (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}).\end{aligned}$$

If  $\boldsymbol{\beta}$  were known, the BLUP of  $\mathbf{b}^T \boldsymbol{\gamma}$  would be

$$\mathbf{h}_1^T (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{b}^T \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}).$$

Since  $\boldsymbol{\beta}$  is unknown, it therefore makes sense to replace it by its BLUE. Thus the BLUP of  $\mathbf{b}^T \boldsymbol{\gamma}$  should be  $\mathbf{b}^T \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})$  and the BLUP of  $\boldsymbol{\gamma}$  should be  $\hat{\boldsymbol{\gamma}} = \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})$  (proof given below). We leave it to the reader to show that the BLUE of  $\boldsymbol{\gamma}$  is unique.

We now write down expressions for the BLUE and BLUP of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  in the following result.

**Theorem 11.10.1.** *The BLUE of  $\boldsymbol{\beta}$  and the BLUP for  $\boldsymbol{\gamma}$  in the mixed linear model are given by*

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y}, \quad \hat{\boldsymbol{\gamma}} = \mathbf{D}^{-1} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}).$$

*Remark 11.10.1.* Unlike in the Gauss-Markov model discussed in Section 11.3 of this chapter, the normal equations for estimating  $\boldsymbol{\beta}$  are no longer of the form  $\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{Y}$ . Except in some balanced random- and mixed-effects models, as an example given below will show, the BLUE of  $\boldsymbol{\beta}$  now depends on  $\boldsymbol{\Sigma}$ , the unknown covariance matrix of  $\mathbf{Y}$ . Similarly prediction of  $\boldsymbol{\gamma}$  also requires the knowledge of  $\boldsymbol{\Sigma}$ . In practice, however,  $\boldsymbol{\Sigma}$  is not known and has to be estimated from the data. In order to obtain the approximate BLUE  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  and the BLUP  $\hat{\boldsymbol{\gamma}}$  of  $\boldsymbol{\gamma}$ , we need to use the estimate  $\hat{\boldsymbol{\Sigma}}$  of  $\boldsymbol{\Sigma}$  in the formulas for  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$ . The problem of estimating  $\boldsymbol{\Sigma}$  which involves estimation of the variance components will be discussed later.

The following result gives necessary and sufficient conditions for  $\hat{\boldsymbol{\gamma}}$  to be BLUP of  $\boldsymbol{\gamma}$ . Its proof is left as an exercise.

**Lemma 11.10.1.** *For the mixed model,  $\hat{\boldsymbol{\gamma}}$  is BLUP for  $\boldsymbol{\gamma}$  if and only if for any  $\mathbf{b} \in \mathbb{R}^m$ ,*

- (i)  $\mathbb{E}[\mathbf{b}^T \hat{\boldsymbol{\gamma}} - \mathbf{b}^T \boldsymbol{\gamma}] = 0$  and
- (ii)  $\text{Cov}[\mathbf{b}^T \hat{\boldsymbol{\gamma}} - \mathbf{b}^T \boldsymbol{\gamma}, \mathbf{l}^T \mathbf{Y}] = 0$  for any  $\mathbf{l} \in \mathbb{R}^n$  satisfying the condition  $\mathbf{X}^T \mathbf{l} = \mathbf{0}$ .

Now let us check if  $\hat{\boldsymbol{\gamma}} = \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})$ , where  $\hat{\boldsymbol{\beta}}$  is the BLUE, is indeed the BLUP of  $\boldsymbol{\gamma}$ . Since  $\hat{\boldsymbol{\beta}}$  is unbiased for  $\boldsymbol{\beta}$ , Condition (i) of Lemma 11.10.1 holds. If Condition (ii) of this lemma also holds, then  $\hat{\boldsymbol{\gamma}}$  is the BLUP of  $\boldsymbol{\gamma}$ . Denoting  $\mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{R} \mathbf{Y}$  where  $\mathbf{R} = \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1}$ , we have  $\hat{\boldsymbol{\gamma}} = \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{I} - \mathbf{R}) \mathbf{Y}$ . Therefore for any  $\mathbf{l} \in \mathbb{R}^n$  with  $\mathbf{X}^T \mathbf{l} = \mathbf{0}$ , we have

$$\begin{aligned}
\text{Cov}[\mathbf{b}^T \hat{\boldsymbol{\gamma}} - \mathbf{b}^T \boldsymbol{\gamma}, \mathbf{l}^T \mathbf{Y}] &= \text{Cov}[\mathbf{b}^T \hat{\boldsymbol{\gamma}}, \mathbf{l}^T \mathbf{Y}] - \text{Cov}[\mathbf{b}^T \boldsymbol{\gamma}, \mathbf{l}^T \mathbf{Y}] \\
&= \text{Cov}[\mathbf{b}^T \hat{\boldsymbol{\gamma}}, \mathbf{l}^T \mathbf{Y}] - \mathbf{b}^T \mathbf{D} \mathbf{Z}^T \mathbf{l} \\
&= \text{Cov}[\mathbf{b}^T \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{I} - \mathbf{X} \mathbf{R}) \mathbf{Y}, \mathbf{l}^T \mathbf{Y}] - \mathbf{b}^T \mathbf{D} \mathbf{Z}^T \mathbf{l} \\
&= \mathbf{b}^T \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{I} - \mathbf{X} \mathbf{R}) \boldsymbol{\Sigma} \mathbf{l} - \mathbf{b}^T \mathbf{D} \mathbf{Z}^T \mathbf{l} = -\mathbf{b}^T \mathbf{D} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} \mathbf{R} \boldsymbol{\Sigma} \mathbf{l}.
\end{aligned}$$

Since  $\mathbf{X}^T \mathbf{l} = \mathbf{0}$ , we have

$$\mathbf{R} \boldsymbol{\Sigma} \mathbf{l} = \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \mathbf{l} = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{l} = \mathbf{0}.$$

This shows that Condition (ii) of [Lemma 11.10.1](#) holds and hence  $\hat{\boldsymbol{\gamma}}$  is the BLUP of  $\boldsymbol{\gamma}$ .

### 11.10.2 Mixed Model Equations

We have already seen that the expressions for the BLUE  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  and the BLUP  $\hat{\boldsymbol{\gamma}}$  of  $\boldsymbol{\gamma}$  involve the inverse of  $\text{Cov}[\mathbf{Y}] = \boldsymbol{\Sigma}$ , which is an  $n \times n$  matrix, where  $n$  is the total number of observations. In many cases,  $\boldsymbol{\Sigma}$  does not have a simple expression, and, if  $n$  is large, calculation of  $\boldsymbol{\Sigma}^{-1}$  may be quite time consuming even on a modern computer. Are there some simpler formulas for obtaining  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$ ? The answer is yes and it turns out that  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$  can be obtained by solving the so-called *Mixed Model Equations*, which require inversions of matrices whose dimensions are much smaller than that of  $\boldsymbol{\Sigma}$ . Let us first see what these equations are, provide intuitive arguments which lead to them, and then show that their solutions are the BLUE and the BLUP.

Mixed model equations take the form (assuming that  $\mathbf{D} = \text{Cov}[\boldsymbol{\gamma}]$  is positive definite)

$$\begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{X} & \sigma_0^2 \mathbf{D}^{-1} + \mathbf{Z}^T \mathbf{Z} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \mathbf{Y} \\ \mathbf{Z}^T \mathbf{Y} \end{pmatrix}. \quad (9)$$

#### Theorem 11.10.2.

- (a) *The mixed model equations have a unique solution.*
- (b) *Let  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\gamma}}$  be the solution to the mixed model equations. Then  $\hat{\boldsymbol{\beta}}$  is the BLUE for  $\boldsymbol{\beta}$  and  $\hat{\boldsymbol{\gamma}}$  is the BLUP for  $\boldsymbol{\gamma}$ .*

The proof of this theorem is postponed to [Section 11.10.3](#). We now give a simple example to show how the mixed model equations can be used in the one-factor random-effects model.

**Example 11.10.6.** Consider a one-factor random-effects model described in [Example 11.10.1](#). Here  $\mathbf{X}$  is an  $n \times 1$  vector of 1's, where  $n = n_1 + \dots + n_k$  is the total number of observations. If  $\pi_i = n_i / (\sigma^2 / \sigma_1^2 + n_i)$ , then we show below that the BLUE and BLUP for  $\mu$  and  $\boldsymbol{\gamma}_i$  are given by (denoting the random effects by  $\{\boldsymbol{\gamma}_i\}$  instead of  $\{\alpha_i\}$ )

$$\hat{\mu} = \sum (1 - \pi_i) n_i \bar{Y}_i / \sum (1 - \pi_i) n_i, \quad \hat{\boldsymbol{\gamma}}_i = \pi_i (\bar{Y}_i - \hat{\mu}).$$

For a balanced model, we have  $\hat{\mu} = \bar{Y}_{..}$  and  $\hat{\gamma}_i = \pi(\bar{Y}_i - \bar{Y}_{..})$  with  $\pi = n_0/(\sigma^2/\sigma_1^2 + n_0)$ , where  $n_0 = n_1 = \dots = n_k$ . In this case, the BLUE of  $\mu$  does not depend on the variance components  $\sigma_1^2$  and  $\sigma^2$ . Also note that for a balanced fixed-effects model with the constraint  $\sum \gamma_i = 0$ , the least squares estimate of  $\gamma_i$  is  $\bar{Y}_i - \bar{Y}_{..}$ . The BLUP of  $\gamma_i$  in the random-effects model is obtained by shrinking  $\bar{Y}_i - \bar{Y}_{..}$  toward zero (since  $0 < \pi < 1$ ). If  $\sigma_1 \rightarrow \infty$ , then  $\pi \rightarrow 1$  and the BLUP is then the same as the BLUE of  $\gamma_i$  for the fixed effects case. In a sense then, the mixed model is the same as the fixed-effects model when  $\sigma_1 = \infty$ .

Let us now see how we can obtain the estimates given above using the mixed model equations. Denote  $Y_{..} = n\bar{Y}_{..}$ , and  $Y_i = n_i\bar{Y}_i$ ,  $i = 1, \dots, k$ . Note that in this case,  $\Sigma = \sigma_1^2 \mathbf{Z}\mathbf{Z}^T + \sigma^2 \mathbf{I}$  with  $\mathbf{D} = \sigma_1^2 \mathbf{I}$ . So we have  $\mathbf{X}^T \mathbf{X} = n\mathbf{I}$ ,  $\mathbf{X}^T \mathbf{Z} = (n_1, \dots, n_k)$ , and  $\mathbf{Z}^T \mathbf{Z}$  is a  $k \times k$  diagonal matrix with diagonal elements  $n_1, \dots, n_k$ . So the mixed model equations are

$$\begin{pmatrix} n & n_1 & \cdots & n_k \\ n_1 & \sigma^2/\sigma_1^2 + n_1 & 0 & 0 \\ \vdots & 0 & \ddots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ n_k & 0 & 0 & \sigma^2/\sigma_1^2 + n_k \end{pmatrix} \begin{pmatrix} \mu \\ \gamma_1 \\ \vdots \\ \gamma_k \end{pmatrix} = \begin{pmatrix} Y_{..} \\ Y_1 \\ \vdots \\ Y_k \end{pmatrix},$$

or

$$\begin{aligned} n\mu + n_1\gamma_1 + \cdots + n_k\gamma_k &= Y_{..}, \text{ and} \\ n_i\mu + (\sigma^2/\sigma_1^2 + n_i)\gamma_i &= Y_i, \quad i = 1, \dots, k. \end{aligned}$$

The last  $k$  equations can be written as

$$\pi_i\mu + \gamma_i = \pi_i(\bar{Y}_i - \mu), \text{ ie, } \mu = \pi_i(\bar{Y}_i - \gamma_i), \quad i = 1, \dots, k.$$

Substitute  $\gamma_i$  by  $\pi_i(\bar{Y}_i - \mu)$  in the first equation and solve for  $\mu$  to obtain

$$\hat{\mu} = \sum(1 - \pi_i)n_i\bar{Y}_i / \sum(1 - \pi_i)n_i \text{ and } \hat{\gamma}_i = \pi_i(\bar{Y}_i - \hat{\mu}).$$

### 11.10.3 Motivation for Mixed Model Equations

We now provide a motivation for the mixed model equations under the assumption of normality. Let us assume that  $\boldsymbol{\gamma} \sim N_q(\mathbf{0}, \mathbf{D})$ ,  $\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$ , and that  $\boldsymbol{\epsilon}$  and  $\boldsymbol{\gamma}$  are independent. The main idea is to treat  $\boldsymbol{\gamma}$  as a parameter even though it is random. Note that conditional on  $\boldsymbol{\gamma}$ ,  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}, \sigma^2 \mathbf{I})$ . The joint pdf of  $\mathbf{Y}$  and  $\boldsymbol{\gamma}$  can be written as

$$\begin{aligned} f_{\mathbf{Y}, \boldsymbol{\gamma}}(\mathbf{y}, \boldsymbol{\gamma}) &= f_{\mathbf{Y}|\boldsymbol{\gamma}}(\mathbf{y}|\boldsymbol{\gamma})f_{\boldsymbol{\gamma}}(\boldsymbol{\gamma}) \\ &= c \exp\left[-\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\gamma}\|^2/(2\sigma^2)\right] \exp\left(-\boldsymbol{\gamma}^T \mathbf{D}^{-1} \boldsymbol{\gamma}/2\right), \end{aligned}$$

where the constant  $c$  depends on  $\sigma$  and  $\mathbf{D}$ . If we treat  $\boldsymbol{\gamma}$  as nonrandom, and try to estimate  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , then the likelihood function is

$$L(\boldsymbol{\beta}, \boldsymbol{\gamma}) = c \exp\left[-\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\gamma}\|^2/(2\sigma^2)\right] \exp\left(-\boldsymbol{\gamma}^T \mathbf{D}^{-1} \boldsymbol{\gamma}/2\right).$$

Maximizing  $L$  with respect to  $\beta$  and  $\gamma$  is equivalent to minimizing  $-2 \log L$  with respect to  $\beta$  and  $\gamma$ . Differentiating  $-2 \log L$  with respect to  $\beta$  and  $\gamma$  and equating the derivatives to zero lead to the following equations

$$\begin{aligned}\partial(-2 \log L) / \partial \beta &= 0, \text{ ie, } X^T X \beta + X^T Z \gamma = X^T Y, \text{ and} \\ \partial(-2 \log L) / \partial \gamma &= 0, \text{ ie, } Z^T X + (Z^T Z + \sigma^2 D^{-1}) \gamma = Z^T Y,\end{aligned}$$

and these are the mixed model equations given in Eq. (9).

The proof of [Theorem 11.10.2](#) requires the following result which can be derived from the Sherman-Morrison formula ([Section B.1](#)). The proof of this lemma is left as an exercise.

**Lemma 11.10.2.**

- (a)  $\Sigma^{-1} = \sigma^{-2} [I - Z(Z^T Z + \sigma^2 D^{-1})^{-1} Z^T]$
- (b)  $\Sigma^{-1} = \sigma^{-2} Z(Z^T Z + \sigma^2 D^{-1})^{-1} D^{-1}$ ,
- (c)  $(Z^T Z + \sigma^2 D^{-1})^{-1} Z^T = D Z^T \Sigma^{-1}$ .

*Proof of Theorem 11.10.2.* Part (a) is not difficult to check and it is left as an exercise.

In order prove part (b), it is enough to show that the solutions of  $\hat{\beta}$  and  $\hat{\gamma}$  of the mixed model equations are the same as those in [Theorem 11.10.1](#).

If  $\hat{\beta}$  and  $\hat{\gamma}$  are the solutions of the mixed model equations, then

$$\begin{aligned}Z^T X \hat{\beta} + (Z^T Z + \sigma^2 D^{-1}) \hat{\gamma} &= Z^T Y, \text{ and hence} \\ \hat{\gamma} &= (Z^T Z + \sigma^2 D^{-1})^{-1} Z^T (Y - X \hat{\beta}) = D Z^T \Sigma^{-1} (Y - X \hat{\beta}),\end{aligned}$$

using part (c) of [Lemma 11.10.2](#). Thus,  $\hat{\gamma}$  has the same form as in [Theorem 11.10.1](#) and hence this would be the BLUP if we can show that the solution  $\hat{\beta}$  of the mixed model equations is indeed the BLUE of  $\beta$ .

The first set of equations in the mixed model equations are

$$X^T X \hat{\beta} + X^T Z \hat{\gamma} = X^T Y.$$

Substituting the expression of  $\hat{\gamma}$  in this equation we have

$$\begin{aligned}X^T X \hat{\beta} + X^T Z D Z^T \Sigma^{-1} (Y - X \hat{\beta}) &= X^T Y, \text{ or} \\ X^T [I - Z D Z^T \Sigma^{-1}] X \hat{\beta} &= X^T [I - Z D Z^T \Sigma^{-1}] Y.\end{aligned}$$

Since

$$I - Z D Z^T \Sigma^{-1} = I - [\Sigma - \sigma^2 I] \Sigma^{-1} = \sigma^2 \Sigma^{-1},$$

we have

$$X^T \sigma^2 \Sigma^{-1} X \hat{\beta} = X^T \sigma^2 \Sigma^{-1} Y, \text{ or } \hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y.$$

The expression is the same as in [Theorem 11.10.1](#). Thus the solution  $\hat{\beta}$  of the mixed model equations is the BLUE of  $\beta$ .  $\square$

#### 11.10.4 Estimation of Variance Components

For the model given in Eq. (7)

$$\mathbf{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}, \text{ and}$$

$$\Sigma = \text{Cov}[\mathbf{Y}] = \sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1^T + \cdots + \sigma_r^2 \mathbf{Z}_r \mathbf{Z}_r^T + \sigma_0^2 I.$$

Note that we have changed the notations a bit and now  $\sigma^2$  is denoted by  $\sigma_0^2$ . We now focus on the problem of estimating the variance components  $\sigma_0^2, \sigma_1^2, \dots, \sigma_r^2$ .

There are a number of well-known methods for estimating the variance components including

- (i) Henderson's method III,
- (ii) Maximum likelihood,
- (iii) Restricted maximum likelihood (REML), and
- (iv) MINQUE (minimum norm quadratic unbiased estimation).

Here we discuss only the first three methods. Computer packages such as R can be used to estimate the variance components using these procedures. Henderson's method and MINQUE do not require any distributional assumptions, whereas the maximum likelihood and the REML methods require the assumptions of normality. It should be pointed out that these methods may not produce the same estimates of the variance components. In unbalanced cases, the estimates obtained by employing any of these procedures usually do not have explicit expressions. Detailed discussion on all these methods can be found in the book "Linear Models" by Searle [60] (Chapters 9–11).

##### *Henderson's Method III*

Suppose that we want to estimate one of the variance components, say  $\sigma_r^2$ . Let  $\mathbf{Q}$  be the projection on the column space of the augmented matrix  $[\mathbf{X}, \mathbf{Z}_1, \dots, \mathbf{Z}_r]$  and let  $\mathbf{Q}_r$  be the projection on the column space of the matrix  $[\mathbf{X}, \mathbf{Z}_1, \dots, \mathbf{Z}_{r-1}]$ . Then  $\mathbf{Q}\mathbf{Y}$  is the vector of fitted values when we fit the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1\boldsymbol{\gamma}_1 + \cdots + \mathbf{Z}_r\boldsymbol{\gamma}_r + \boldsymbol{\epsilon}$  pretending that  $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_r$  are nonrandom unknown parameters. Similarly  $\mathbf{Q}_r\mathbf{Y}$  is the vector of fitted values when we fit the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1\boldsymbol{\gamma}_1 + \cdots + \mathbf{Z}_{r-1}\boldsymbol{\gamma}_{r-1} + \boldsymbol{\epsilon}$  pretending that  $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_{r-1}$  are nonrandom parameters. Note that an unbiased estimate of  $\sigma_0^2$  is given by  $\hat{\sigma}_0^2 = \|\mathbf{Y} - \mathbf{Q}\mathbf{Y}\|^2/(n - \text{rank}(\mathbf{Q}))$ . Henderson's estimate of  $\sigma_r^2$  is given by

$$\hat{\sigma}_r^2 = \left[ \|\mathbf{Q}\mathbf{Y} - \mathbf{Q}_r\mathbf{Y}\|^2 - \text{trace}(\mathbf{Q} - \mathbf{Q}_r)\hat{\sigma}_0^2 \right] / \text{trace}(\mathbf{Z}_r^T(\mathbf{I} - \mathbf{Q}_r)\mathbf{Z}_r). \quad (10)$$

Since  $\sigma_r^2 > 0$ , its estimate is usually taken to be  $\max(\hat{\sigma}_r^2, 0)$ .

Let us now see how this estimate is derived. First of all note that  $(\mathbf{Q} - \mathbf{Q}_r)\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$  and  $(\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_i = \mathbf{0}$ ,  $i = 1, \dots, r-1$ . Hence

$$\begin{aligned} (\mathbf{Q} - \mathbf{Q}_r)\mathbf{Y} &= (\mathbf{Q} - \mathbf{Q}_r)(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1\boldsymbol{\gamma}_1 + \cdots + \mathbf{Z}_r\boldsymbol{\gamma}_r + \boldsymbol{\epsilon}) \\ &= (\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\boldsymbol{\gamma}_r + (\mathbf{Q} - \mathbf{Q}_r)\boldsymbol{\epsilon}. \end{aligned}$$

Then

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{QY} - \mathbf{Q}_r\mathbf{Y}\|^2\right] &= \mathbb{E}\left[\|(\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\boldsymbol{\gamma}_r + (\mathbf{Q} - \mathbf{Q}_r)\boldsymbol{\varepsilon}\|^2\right] \\ &= \mathbb{E}\left[\|(\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\boldsymbol{\gamma}_r\|^2\right] + \mathbb{E}\left[\|(\mathbf{Q} - \mathbf{Q}_r)\boldsymbol{\varepsilon}\|^2\right] \\ &= \sigma_r^2 \text{trace}\left((\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\mathbf{Z}_r^T(\mathbf{Q} - \mathbf{Q}_r)^T\right) + \sigma_0^2 \text{trace}\left((\mathbf{Q} - \mathbf{Q}_r)(\mathbf{Q} - \mathbf{Q}_r)^T\right) \\ &= \sigma_r^2 \text{trace}\left(\mathbf{Z}_r^T(\mathbf{Q} - \mathbf{Q}_r)^T(\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\right) + \sigma_0^2 \text{trace}\left((\mathbf{Q} - \mathbf{Q}_r)(\mathbf{Q} - \mathbf{Q}_r)^T\right). \end{aligned}$$

Now use the facts that  $\mathbf{Q}$  and  $\mathbf{Q}_r$  are symmetric and  $\mathbf{Q} - \mathbf{Q}_r$  is also a projection, and consequently,  $(\mathbf{Q} - \mathbf{Q}_r)^T(\mathbf{Q} - \mathbf{Q}_r) = \mathbf{Q} - \mathbf{Q}_r$  and  $(\mathbf{Q} - \mathbf{Q}_r)(\mathbf{Q} - \mathbf{Q}_r)^T = \mathbf{Q} - \mathbf{Q}_r$ . So we have

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{QY} - \mathbf{Q}_r\mathbf{Y}\|^2\right] &= \sigma_r^2 \text{trace}\left(\mathbf{Z}_r^T(\mathbf{Q} - \mathbf{Q}_r)\mathbf{Z}_r\right) + \sigma_0^2 \text{trace}(\mathbf{Q} - \mathbf{Q}_r) \\ &= \sigma_r^2 \text{trace}\left(\mathbf{Z}_r^T(I - \mathbf{Q}_r)\mathbf{Z}_r\right) + \sigma_0^2 \text{trace}(\mathbf{Q} - \mathbf{Q}_r). \end{aligned}$$

The last step is justified since  $\mathbf{Z}_r^T\mathbf{QZ}_r = \mathbf{Z}_r^T\mathbf{Q}^T\mathbf{QZ}_r = \mathbf{Z}_r^T\mathbf{Z}_r$ . From the last expression we now see that an unbiased estimate of  $\sigma_r^2$  is of the form given above.

Except in the case of balanced models, Henderson's estimates for the variance components do not have nice forms. Here we give an example for the one-factor case.

**Example 11.10.7.** Consider a one-factor (random) ANOVA model as in Example 11.10.6.

Here

$$\begin{aligned} \|\mathbf{Y} - \mathbf{QY}\|^2 &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{i\cdot})^2 = SSE, \text{ and} \\ \|\mathbf{QY} - \mathbf{Q}_1\mathbf{Y}\|^2 &= \sum_i \sum_j (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2 = \sum n_i (\bar{Y}_{i\cdot} - \bar{Y}_{..})^2 = SSTR. \end{aligned}$$

Assuming that  $n = n_1 + \dots + n_k$  is the total number of observations, we have  $n - \text{rank}(\mathbf{Q}) = n - k$  and  $\text{trace}(\mathbf{Q} - \mathbf{Q}_1) = \text{rank}(\mathbf{Q}) - \text{rank}(\mathbf{Q}_1) = k - 1$ . So an unbiased estimate of  $\sigma_0^2$  is given by

$$\hat{\sigma}_0^2 = \|\mathbf{Y} - \mathbf{QY}\|^2 / (n - \text{rank}(\mathbf{Q})) = SSE / (n - k) = MSE.$$

Now let  $\bar{\gamma} = \sum n_i \gamma_i / n$ . Then

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{QY} - \mathbf{Q}_1\mathbf{Y}\|^2\right] &= \mathbb{E}[SSTR] \\ &= \mathbb{E}\left[\sum n_i (\gamma_i - \bar{\gamma} + \bar{\varepsilon}_{i\cdot} - \varepsilon_{..})^2\right] \\ &= (n - \sum n_i^2 / n) \sigma_1^2 + (k - 1) \sigma_0^2. \end{aligned}$$

So an unbiased estimate of  $\sigma_1^2$  is given by

$$\hat{\sigma}_1^2 = [SSTR - (k - 1)\hat{\sigma}_0^2] / [n - \sum n_i^2 / n]$$

$$\begin{aligned}
&= [SSTR - (k-1)MSE] / \left[ n - \sum n_i^2/n \right] \\
&= (k-1)[MSTR - MSE] / \left[ n - \sum n_i^2/n \right].
\end{aligned}$$

In the balanced case,  $n = n_0a$  where  $n_i \equiv n_0$  for all  $i$ . It is easy to check that the Henderson's estimate of  $\sigma_1^2$  is then given by  $\hat{\sigma}_1^2 = [MSTR - MSE]/n_0$ .

## Maximum Likelihood

The maximum likelihood method (under the assumption of joint normality of  $\gamma_1, \dots, \gamma_r$  and  $\boldsymbol{\varepsilon}$ ) jointly estimates  $\boldsymbol{\beta}$  and the variance components. We rewrite the mixed linear model given in Eq. (7) as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \sum_{i=0}^r \mathbf{Z}_i \boldsymbol{\gamma}_i, \quad (11a)$$

$$\mathbb{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}, \text{ and } \boldsymbol{\Sigma}(\boldsymbol{\theta}) = \text{Cov}[\mathbf{Y}] = \sum_{i=0}^r \sigma_i^2 \mathbf{Z}_i \mathbf{Z}_i^T, \quad (11b)$$

where  $\mathbf{Z}_0 = \mathbf{I}$ ,  $\boldsymbol{\gamma}_0 = \boldsymbol{\varepsilon}$ , and  $\boldsymbol{\theta} = (\sigma_0^2, \sigma_1^2, \dots, \sigma_r^2)^T$ . So we have  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$ . The likelihood therefore is

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}) = c \left[ 1/|\boldsymbol{\Sigma}(\boldsymbol{\theta})|^{1/2} \right] \exp \left[ -(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})/2 \right],$$

where  $|\boldsymbol{\Sigma}(\boldsymbol{\theta})| = \det[\boldsymbol{\Sigma}(\boldsymbol{\theta})]$  and  $c > 0$  is a constant that does not depend on  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$ . Maximizing the likelihood with respect to  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$  is equivalent to minimizing  $-2 \log L$  with respect to  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$ . Note that

$$-2 \log L(\boldsymbol{\beta}, \boldsymbol{\theta}) = -2 \log(c) + (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \log(|\boldsymbol{\Sigma}(\boldsymbol{\theta})|).$$

We need to differentiate  $-2 \log L$  with respect to  $\boldsymbol{\beta}$  and  $\theta_i$ ,  $i = 0, \dots, r$ , and equate the derivatives to zero. Calculations will show that

$$\begin{aligned}
\partial(-2 \log L)/\partial \boldsymbol{\beta} &= 0, \text{ ie, } \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Y}, \text{ and} \\
\partial(-2 \log L)/\partial \theta_i &= 0, \text{ ie,} \\
(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[ \partial \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} / \partial \theta_i \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \partial(\log |\boldsymbol{\Sigma}(\boldsymbol{\theta})|) / \partial \theta_i &= 0, \quad i = 0, \dots, r.
\end{aligned}$$

Note that the first set of equations are the same as the normal equations (for estimating  $\boldsymbol{\beta}$ ) described in Section 11.10.1. The second set of equations need simplifications. Noting that  $\partial \boldsymbol{\Sigma}(\boldsymbol{\theta})/\partial \theta_i = \mathbf{Z}_i \mathbf{Z}_i^T$  and using results from Section B.5, we have

$$\begin{aligned}
\partial \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} / \partial \theta_i &= -\boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} [\partial \boldsymbol{\Sigma}(\boldsymbol{\theta}) / \partial \theta_i] \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \\
&= -\boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i \mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1}, \text{ and} \\
\partial(\log |\boldsymbol{\Sigma}(\boldsymbol{\theta})|) / \partial \theta_i &= \text{trace}(\boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \partial \boldsymbol{\Sigma}(\boldsymbol{\theta}) / \partial \theta_i) \\
&= \text{trace}(\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i).
\end{aligned}$$

Thus the likelihood equation involving derivative with respect to  $\theta_i$  turns out to be

$$\begin{aligned} 0 &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[ \partial \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} / \partial \theta_i \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \partial(\log |\boldsymbol{\Sigma}(\boldsymbol{\theta})|) / \partial \theta_i \\ &= -(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i \mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \text{trace}(\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i) \\ &= -\|\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2 + \text{trace}(\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i). \end{aligned}$$

So the likelihood equations are

$$\mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Y}, \quad (12a)$$

$$\text{trace}(\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_i) = \|\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2, \quad i = 0, \dots, r. \quad (12b)$$

These equations have no explicit solutions except in some balanced cases and iterative methods are used to solve them numerically.

### *Restricted Maximum Likelihood*

REML is a variant of the maximum likelihood method whereby the issue of estimation of  $\boldsymbol{\beta}$  is entirely bypassed and the focus is entirely on estimating the variance components. In this method, the estimate of  $\boldsymbol{\theta} = (\sigma_0^2, \dots, \sigma_r^2)^T$  is obtained by solving the equations

$$\text{trace}(\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{I} - \mathbf{M}(\boldsymbol{\theta})) \mathbf{Z}_i) = \|\mathbf{Z}_i^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{I} - \mathbf{M}(\boldsymbol{\theta})) \mathbf{Y}\|^2, \quad i = 0, \dots, r, \quad (13)$$

where  $\mathbf{M}(\boldsymbol{\theta}) = \mathbf{X} \left( \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1}$ . These equations usually have no explicit solutions and iterative methods are employed in numerical computations. Once an estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  is available,  $\boldsymbol{\beta}$  can be estimated by solving the approximate normal equations

$$\mathbf{X}^T \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^T \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{Y}.$$

We will provide a justification for the REML equations. Since  $\text{rank}(\mathbf{X}) = p$ , we can find a matrix  $\mathbf{B}$  of order  $(n - p) \times n$  which has rank  $n - p$  and it satisfies the equation  $\mathbf{BX} = \mathbf{0}$ . Let  $\tilde{\mathbf{Y}} = \mathbf{BY}$ ,  $\tilde{\mathbf{Z}}_i = \mathbf{BZ}_i$ ,  $i = 0, \dots, r$ , premultiplying both sides of Eq. (11a) leads to a modified model

$$\begin{aligned} \tilde{\mathbf{Y}} &= \sum_{i=0}^r \tilde{\mathbf{Z}}_i \boldsymbol{\gamma}_i, \text{ and} \\ \text{Cov}[\tilde{\mathbf{Y}}] &= \sum_{i=0}^r \sigma_i^2 \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^T = \tilde{\boldsymbol{\Sigma}}(\boldsymbol{\theta}), \text{ say.} \end{aligned}$$

Since  $\tilde{\mathbf{Y}} \sim N_{n-p}(\mathbf{0}, \tilde{\boldsymbol{\Sigma}}(\boldsymbol{\theta}))$ , the likelihood (based on  $\tilde{\mathbf{Y}}$ ) is

$$L(\boldsymbol{\theta}) = c \left[ 1/|\tilde{\boldsymbol{\Sigma}}(\boldsymbol{\theta})|^{1/2} \right] \exp \left[ -\tilde{\mathbf{Y}}^T \tilde{\boldsymbol{\Sigma}}(\boldsymbol{\theta})^{-1} \tilde{\mathbf{Y}} / 2 \right],$$

where the constant  $c > 0$  does not depend on  $\boldsymbol{\theta}$ . In order to obtain the MLE for  $\boldsymbol{\theta}$  in this case, we differentiate  $-2 \log L$  with respect to  $\theta_i$ ,  $i = 0, \dots, r$ , and equate the derivatives to

zero. Then, as in Eq. (12b), we have

$$\text{trace}(\tilde{\mathbf{Z}}_i^T \tilde{\Sigma}(\theta)^{-1} \tilde{\mathbf{Z}}_i) = \|\tilde{\mathbf{Z}}_i^T \tilde{\Sigma}(\theta)^{-1} \tilde{\mathbf{Y}}\|^2, \quad i = 0, \dots, r.$$

Since  $\tilde{\mathbf{Z}}_i = \mathbf{B}\mathbf{Z}_i$  and  $\tilde{\Sigma}(\theta) = \mathbf{B}\Sigma(\theta)\mathbf{B}^T$ , the  $i$ th equation is

$$\text{trace}(\mathbf{Z}_i^T \mathbf{B}^T (\mathbf{B}\Sigma(\theta)\mathbf{B}^T)^{-1} \mathbf{B}\mathbf{Z}_i) = \|\mathbf{Z}_i^T \mathbf{B}^T (\mathbf{B}\Sigma(\theta)\mathbf{B}^T)^{-1} \mathbf{B}\mathbf{Y}\|^2.$$

The important fact is that the matrix  $\mathbf{B}^T (\mathbf{B}\Sigma(\theta)\mathbf{B}^T)^{-1} \mathbf{B}$  does not depend on the choice  $\mathbf{B}$  as long as  $\mathbf{B}\mathbf{X} = \mathbf{0}$  and  $\text{rank}(\mathbf{B}) = n - \text{rank}(\mathbf{X})$ . The following result turns out to be true.

**Lemma 11.10.3.** *Assume that  $\mathbf{B}\mathbf{X} = \mathbf{0}$  and  $\text{rank}(\mathbf{B}) = n - \text{rank}(\mathbf{X})$ . Let  $\mathbf{M}(\theta) = \mathbf{X}(\mathbf{X}^T \Sigma(\theta)^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma(\theta)^{-1}$ . Then*

$$\mathbf{B}^T (\mathbf{B}\Sigma(\theta)\mathbf{B}^T)^{-1} \mathbf{B} = \Sigma(\theta)^{-1} (\mathbf{I} - \mathbf{M}(\theta)).$$

Using Lemma 11.10.3, we can show that the REML equations are

$$\text{trace}(\mathbf{Z}_i^T \Sigma(\theta)^{-1} (\mathbf{I} - \mathbf{M}(\theta)) \mathbf{Z}_i) = \|\mathbf{Z}_i^T \Sigma(\theta)^{-1} (\mathbf{I} - \mathbf{M}(\theta)) \mathbf{Y}\|^2, \quad i = 1, \dots, r.$$

*Proof of Lemma 11.10.3.* Simplifying the notations, writing  $\Sigma$  instead of  $\Sigma(\theta)$  and denoting  $\mathbf{R} = \Sigma^{1/2} \mathbf{B}^T$ , where  $\Sigma^{1/2}$  is a symmetric square root of  $\Sigma$ , we have

$$\mathbf{B}^T (\mathbf{B}\Sigma\mathbf{B}^T)^{-1} \mathbf{B} = \Sigma^{-1/2} \mathbf{R} (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \Sigma^{-1/2}.$$

Now note that  $\mathbf{R}(\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T$  is a projection matrix. Where does this matrix project onto? Since  $\mathbf{B}\mathbf{X} = \mathbf{0}$  and  $\text{rank}(\mathbf{B}) = n - \text{rank}(\mathbf{X})$ ,  $\mathcal{M}(\mathbf{R})$  is the same as the orthogonal complement of the column space of the matrix  $\Sigma^{-1/2} \mathbf{X}$ . Now the expression of the matrix that projects on the column space of  $\Sigma^{-1/2} \mathbf{X}$  is given by  $\Sigma^{-1/2} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1/2}$ . So  $\mathbf{I} - \Sigma^{-1/2} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1/2}$  projects on the orthogonal complement of  $\Sigma^{-1/2} \mathbf{X}$ . Hence we must have

$$\mathbf{R} (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T = \mathbf{I} - \Sigma^{-1/2} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1/2}.$$

Consequently,

$$\begin{aligned} \mathbf{B}^T (\mathbf{B}\Sigma\mathbf{B}^T)^{-1} \mathbf{B} &= \Sigma^{-1/2} \mathbf{R} (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \Sigma^{-1/2} \\ &= \Sigma^{-1/2} \left( \mathbf{I} - \Sigma^{-1/2} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1/2} \right) \Sigma^{-1/2} \\ &= \Sigma^{-1} (\mathbf{I} - \mathbf{M}), \end{aligned}$$

where  $\mathbf{M} = \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1}$ . This completes the proof.  $\square$

### 11.11 Inference: Examples From Mixed Models

Here we present estimation methods in a few balanced models since the expressions for the estimates in the unbalanced cases involve complicated and cumbersome notations. Except in [Example 11.11.3](#), where the maximum likelihood method is considered, Henderson's method is used throughout to obtain the estimates of the variance components.

**Example 11.11.1** (One-Factor Random Effects). We consider one-factor balanced ANOVA random-effects models as in [Examples 11.10.1](#), [11.10.6](#), and [11.10.7](#), except now we assume that  $\{\gamma_i\}$  are iid  $N(0, \sigma_1^2)$  and  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma_0^2)$ . We have already see in [Examples 11.10.6](#) and [11.10.7](#) that the BLUE for  $\mu$  is  $\hat{\mu} = \bar{Y}_{..}$ ,

$$\begin{aligned} E[SSE] &= (n - k)\sigma_0^2, \text{ where } n = n_0k, \\ E[SSTR] &= (k - 1)n_0\sigma_1^2 + (k - 1)\sigma_0^2, \\ E[MSE] &= \sigma_0^2, E[MSTR] = n_0\sigma_1^2 + \sigma_0^2, \end{aligned}$$

and unbiased estimate of  $\sigma_0^2$  and  $\sigma_1^2$  are

$$\hat{\sigma}_0^2 = MSE \text{ and } \hat{\sigma}_1^2 = (MSTR - MSE)/n_0.$$

Direct calculation will show that  $\hat{\mu} = \bar{Y}_{..} \sim N(\mu, (n_0\sigma_1^2 + \sigma_0^2)/n)$ . So an estimate of  $\text{Var}[\hat{\mu}]$  is then given by  $s^2(\hat{\mu}) = MSTR/n$ . Note that  $SSTR/(n_0\sigma_1^2 + \sigma_0^2) \sim \chi_{k-1}^2$  and  $(\hat{\mu} - \mu)/\sqrt{(n_0\sigma_1^2 + \sigma_0^2)/n} \sim N(0, 1)$ . Since  $\hat{\mu} = \bar{Y}_{..}$  is independent of  $SSTR$  and hence of  $MSTR$ ,

$$(\hat{\mu} - \mu)/s(\hat{\mu}) = \frac{(\hat{\mu} - \mu)/\sqrt{(n_0\sigma_1^2 + \sigma_0^2)/n}}{\sqrt{SSTR/\left[(n_0\sigma_1^2 + \sigma_0^2)(k - 1)\right]}} \sim t_{k-1}$$

and this fact now can be used to construct a confidence interval for  $\mu$ . Unlike in the fixed-effect case, MSTR is being used to estimate  $\text{Var}[\hat{\mu}]$  and, consequently, the df for the  $t$ -distribution is now  $k - 1$  instead of  $n - k$ .

In some cases one may want to test if  $H_0: \sigma_1^2 = 0$  vs  $H_1: \sigma_1^2 \neq 0$ . The  $F$ -statistic for this is  $F = MSTR/MSE$  and  $F \sim F_{k-1, n-k}$  under  $H_0$ .

**Example 11.11.2.** Consider the one-factor random-effects model as in the last example and we want to construct a confidence interval for the intraclass correlation coefficient  $\rho = \sigma_1^2/(\sigma_1^2 + \sigma_0^2)$ , which is also the proportion of variability in the response explained by the random factor. Clearly an estimate of  $\rho$  is given by  $\hat{\rho} = \hat{\sigma}_1^2/(\hat{\sigma}_1^2 + \hat{\sigma}_0^2)$ , where the expressions for  $\hat{\sigma}_1^2$  and  $\hat{\sigma}_0^2$  are as given in the last example. Let us denote  $MSTR/MSE$  by  $F^*$ . A confidence interval for  $\rho$  with confidence coefficient  $1 - \alpha$  is given by  $[L^*, U^*]$ , where

$L^* = L/(1+L)$ ,  $U^* = U/(1+U)$ , and

$$L = (1/n_0) \left[ \frac{F^*}{F_{k-1,n-k,\alpha/2}} - 1 \right] \text{ and } U = (1/n_0) \left[ \frac{F^*}{F_{k-1,n-k,1-\alpha/2}} - 1 \right].$$

Let us see why this is indeed a confidence interval of  $\rho$  with confidence coefficient  $1 - \alpha$ . Note that  $SSTR/(n_0\sigma_1^2 + \sigma_0^2) \sim \chi_{k-1}^2$ ,  $SSE/\sigma_0^2 \sim \chi_{n-k}^2$ , and that  $SSTR$  and  $SSE$  are independent. Hence the random variable

$$F = \frac{\left[ SSTR / (n_0\sigma_1^2 + \sigma_0^2) \right] / (k-1)}{\left[ SSE / \sigma_0^2 \right] / (n-k)} = \left[ \sigma_0^2 / (n_0\sigma_1^2 + \sigma_0^2) \right] F^*$$

has an  $F$ -distribution with  $df = (k-1, n-k)$ . Denote  $c_1 = F_{k-1,n-k,1-\alpha/2}$  and  $c_2 = F_{k-1,n-k,\alpha/2}$ . Then

$$\begin{aligned} 1 - \alpha &= P[c_1 \leq F \leq c_2] \\ &= P[c_1/F^* \leq \sigma_0^2 / (n_0\sigma_1^2 + \sigma_0^2) \leq c_2/F^*] \\ &= P[F^*/c_2 \leq n_0\sigma_1^2 / \sigma_0^2 + 1 \leq F^*/c_1] \\ &= P[(F^*/c_2 - 1)/n_0 \leq \sigma_1^2 / \sigma_0^2 \leq (F^*/c_1 - 1)/n_0] \\ &= P[L \leq \sigma_1^2 / \sigma_0^2 \leq U] \\ &= P[L/(1+L) \leq \sigma_1^2 / (\sigma_1^2 + \sigma_0^2) \leq U/(1+U)] \\ &= P[L^* \leq \rho \leq U^*]. \end{aligned}$$

**Example 11.11.3** (One-Factor Random-Effects Model). The setup here is the same as in the last two examples, but the estimates of the variance components are obtained using the maximum likelihood method. Here  $\mathbf{X}$  is  $n$ -dim vector of 1's,  $\beta = \mu$  is a scalar,  $\mathbf{Z}_0 = \mathbf{I}$ , and

the matrix  $\mathbf{Z}_1 = \begin{pmatrix} \mathbf{1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1} \end{pmatrix}$  is  $n \times k$ , where  $\mathbf{1}$  is the  $n_0$ -dim vector of 1's. Since  $\boldsymbol{\Sigma}(\boldsymbol{\theta}) =$

$\sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1^T + \sigma_0^2 \mathbf{I}$ , it can be checked using the Sherman-Morrison formula (Section B.1) that

$$\boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} = \sigma_0^{-2} [\mathbf{I} - \pi n_0^{-1} \mathbf{Z}_1 \mathbf{Z}_1^T], \text{ where } \pi = n_0 \sigma_1^2 / (n_0 \sigma_1^2 + \sigma_0^2),$$

$$\mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{X} = \sigma_0^{-2} n(1 - \pi), \text{ and}$$

$$\mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\beta) = \sigma_0^{-2} (1 - \pi)(Y.. - n\mu).$$

Thus the first likelihood Eq. (12a)

$$\begin{aligned} \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{X} &= \mathbf{X}^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\beta), \text{ ie,} \\ \sigma_0^{-2} (1 - \pi)(Y.. - n\mu) &= \sigma_0^{-2} n(1 - \pi) \end{aligned}$$

leads to the usual estimate of  $\mu$  (ie,  $\hat{\mu} = \bar{Y}_{..}$ ). In order to obtain the maximum likelihood estimates of  $\sigma_0^2$  and  $\sigma_1^2$ , more calculations are needed. The following can be checked (left as an exercise)

$$\begin{aligned} \text{trace}(\mathbf{Z}_0^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_0) &= \sigma_0^{-2} n(1 - \pi n_0^{-1}), \\ \text{trace}(\mathbf{Z}_1^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} \mathbf{Z}_1) &= \sigma_0^{-2} n(1 - \pi), \\ \|\mathbf{Z}_0^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta})\|^2 &= \sigma_0^{-4} \sum_i \sum_j [Y_{ij} - \bar{Y}_{..} - \pi(\bar{Y}_{i..} - \bar{Y}_{..})]^2 \\ &= \sigma_0^{-4} [SSE + (1 - \pi)^2 SSTR], \text{ and} \\ \|\mathbf{Z}_1^T \boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta})\|^2 &= \sigma_0^{-4} n_0 (1 - \pi)^2 SSTR. \end{aligned}$$

Using the above-simplified expressions, the likelihood Eq. (12b) in this case are

$$\begin{aligned} \sigma_0^{-2} n(1 - \pi n_0^{-1}) &= \sigma_0^{-4} [SSE + (1 - \pi)^2 SSTR], \text{ and} \\ \sigma_0^{-2} n(1 - \pi) &= \sigma_0^{-4} n_0 (1 - \pi)^2 SSTR. \end{aligned}$$

Solutions to these equations lead to the maximum likelihood estimates of the variance components

$$\hat{\sigma}_0^2 = MSE \quad \text{and} \quad \hat{\sigma}_1^2 = \frac{SSTR/k - MSE}{n_0} = \frac{[(k-1)/k]MSTR - MSE}{n_0}.$$

This estimate of  $\sigma_1^2$  is slightly different from Henderson's estimate given in Example 11.10.7.

**Example 11.11.4** (Two-Factor ANOVA, Both Factors Random). The setup here is the same as in Example 11.10.2, but we now assume that it is a balanced ANOVA (ie,  $n_{ij} = n_0$  for all  $i$  and  $j$ ). It can be checked that the BLUE of  $\mu$  here is  $\hat{\mu} = \bar{Y}_{...}$ . Recall that the sums of squares are (Example 11.6.4)

$$\begin{aligned} SSA &= n_0 b \sum_i (\bar{Y}_{i..} - \bar{Y}_{...})^2, \quad SSB = n_0 a \sum_j (\bar{Y}_{.j} - \bar{Y}_{...})^2, \\ SSAB &= n_0 \sum_i \sum_j (\bar{Y}_{ij..} - \bar{Y}_{i..} - \bar{Y}_{.j} + \bar{Y}_{...})^2, \text{ and} \\ SSE &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij..})^2. \end{aligned}$$

It can be shown that

$$\begin{aligned} E[MSA] &= n_0 b \sigma_1^2 + n_0 \sigma_3^2 + \sigma_0^2, \\ E[MSB] &= n_0 a \sigma_2^2 + n_0 \sigma_3^2 + \sigma_0^2, \\ E[MSAB] &= n_0 \sigma_3^2 + \sigma_0^2, \quad \text{and} \quad E[MSE] = \sigma_0^2. \end{aligned}$$

So the unbiased estimates of the variance components are

$$\begin{aligned}\hat{\sigma}_0^2 &= MSE, \quad \hat{\sigma}_3^2 = (MSAB - MSE)/n_0, \\ \hat{\sigma}_1^2 &= (MSA - MSAB)/(n_0 b), \quad \text{and} \quad \hat{\sigma}_2^2 = (MSB - MSAB)/(n_0 a).\end{aligned}$$

Note that  $\hat{\mu} = \bar{Y}_{...}$  is unbiased for  $\mu$  and

$$\text{Var}[\hat{\mu}] = (n_0 b \sigma_1^2 + n_0 a \sigma_2^2 + n_0 \sigma_3^2 + \sigma_0^2)/n,$$

where  $n = n_0 ab$ . So an unbiased estimate of  $\text{Var}[\hat{\mu}]$  is given by  $s^2(\hat{\mu}) = (MSA + MSB - MSAB)/n$ .

**Example 11.11.5** (Two-Factor ANOVA: Factor A Fixed, Factor B Random). We now consider a two-factor balanced mixed-effects model as in [Example 11.10.3](#).

For this model,

$$E[Y_{ijk}] = \mu + \alpha_i \quad \text{and} \quad \text{Var}[Y_{ijk}] = \sigma_1^2 + \sigma_2^2 + \sigma_0^2.$$

It can be shown that the expected values of the mean squares are

$$\begin{aligned}E[MSA] &= n_0 b \sum (\alpha_i - \bar{\alpha})^2 / (a - 1) + n_0 \sigma_2^2 + \sigma_0^2, \\ E[MSB] &= n_0 a \sigma_1^2 + n_0 \sigma_2^2 + \sigma_0^2, \\ E[MSAB] &= n_0 \sigma_2^2 + \sigma_0^2, \quad \text{and} \quad E[MSE] = \sigma_0^2.\end{aligned}$$

So the unbiased estimates of the variance components are

$$\begin{aligned}\hat{\sigma}_0^2 &= MSE, \quad \hat{\sigma}_2^2 = (MSAB - MSE)/n_0, \quad \text{and} \\ \hat{\sigma}_1^2 &= (MSB - MSAB)/(n_0 a).\end{aligned}$$

If we want to test the null hypothesis of no effect of factor A (ie,  $H_0: \alpha_1 = \dots = \alpha_a$  vs  $H_1: \text{not all } \alpha_i's \text{ are equal}$ ), then the test statistic is  $F^* = MSA/MSAB$  with  $df = (a - 1, (a - 1)(b - 1))$ .

For this model, the BLUE of  $\mu_i = \mu + \alpha_i$  is given by  $\bar{Y}_{i..}$ . It can be shown that  $E[\bar{Y}_{i..}] = \mu + \alpha_i$  and

$$\text{Var}[\bar{Y}_{i..}] = (n_0 \sigma_1^2 + n_0 \sigma_2^2 + \sigma_0^2)/(n_0 b).$$

So an unbiased estimate of  $\text{Var}[\bar{Y}_{i..}]$  is given by

$$s^2(\bar{Y}_{i..}) = [(MSB - MSAB)/a + MSAB]/(n_0 b).$$

Note that

$$(MSB - MSAB)/a + MSAB = (1/a)MSB + (1 - 1/a)MSAB.$$

The last quantity is nonnegative since it is a weighted average of  $MSB$  and  $MSAB$ . Also note that

$$\begin{aligned} (MSB - MSAB)/a + MSAB &= (1/a)MSB + (1 - 1/a)MSAB \\ &= [(b-1)MSB + (a-1)(b-1)MSAB]/[a(b-1)] \\ &= [SSB + SSAB]/[a(b-1)] \\ &= n_0 \sum \sum (\bar{Y}_{ij} - \bar{Y}_{i..})^2 / [a(b-1)]. \end{aligned}$$

We denote the sum of squares  $n_0 \sum_i \sum_j (\bar{Y}_{ij} - \bar{Y}_{i..})^2$  by  $SSB(A)$  (nested ANOVA case in Example 11.6.6). This sum of squares has  $df = a(b-1)$  and we denote  $SSB(A)/[a(b-1)]$  by  $MSB(A)$ . Hence we have

$$s^2(\bar{Y}_{i..}) = MSB(A)/(n_0 b).$$

Since  $\bar{Y}_{i..}$  is independent of  $SSB(A)$  and hence of  $MSB(A)$ ,  $(\bar{Y}_{i..} - \mu_i)/s(\bar{Y}_{i..}) \sim t_{a(b-1)}$  and this fact can be used to construct a confidence interval for  $\mu_i$ .

If we want to estimate a contrast  $\theta = \sum c_i \alpha_i$  of  $\{\alpha_i\}$ , then the BLUE is  $\hat{\theta} = \sum c_i \bar{Y}_{i..}$ . It can be shown that

$$E[\hat{\theta}] = \theta \text{ and } \text{Var}[\hat{\theta}] = \left( \sum c_i^2 \right) (n_0 \sigma_2^2 + \sigma_0^2) / (n_0 b).$$

An unbiased estimate of  $\text{Var}[\hat{\theta}]$  is

$$s^2(\hat{\theta}) = \left( \sum c_i^2 \right) MSAB / (n_0 b),$$

and since  $\{\bar{Y}_{i..}, i = 1, \dots, a\}$  are independent of  $SSAB$  and hence of  $MSAB$ ,  $\hat{\theta}$  is also independent of  $MSAB$ . Consequently,  $(\hat{\theta} - \theta) / s(\hat{\theta}) \sim t_{(a-1)(b-1)}$ , and this fact can be used to construct a confidence interval for  $\theta$ .

*Remark 11.11.1.* An alternative modeling scheme for the interactions in the last example assumes that

- (i) for any  $j$ ,  $\sum_{i=1}^a (\alpha\beta)_{ij} = 0$ ,
- (ii) for any  $i$ ,  $(\alpha\beta)_{i1}, \dots, (\alpha\beta)_{ib}$  are iid  $N(0, \sigma_2^2)$ ,
- (iii)  $\text{Cov}[(\alpha\beta)_{ij}, (\alpha\beta)_{i'j}] = -\sigma_2^2/(a-1)$ , for any  $i \neq i', j = 1, \dots, b$ .

In this framework  $\{(\alpha\beta)_{ij}\}$  are no longer iid and the estimates of the variance components may be different from what are given above.

## Exercises

- 11.1.** Consider a one-factor study with  $k$  levels as in [Example 11.2.3](#).
- (a) Express the ANOVA model in a regression setting by creating  $k - 1$  indicator variables for the factor levels. Relate the parameters of this regression model to those of the ANOVA model.
  - (b) Obtain the parameter estimates of the regression model in part (a) and their standard errors.
- 11.2.** Verify the expression of  $E[MSTR]$  given in [Example 11.4.2](#).
- 11.3.** Consider the two-factor ANOVA model given in [Example 11.4.4](#).
- (a) Verify the expressions for  $E[SSA]$ ,  $E[SSB]$ , and  $E[SSAB]$  as given in [Example 11.4.4](#).
  - (b) Check that  $E[MSAB] = \sigma^2$  if and only if  $(\alpha\beta)_{ij} = 0$  for all  $i$  and  $j$ . Similarly, check that  $E[MSA] = \sigma^2$  if and only if  $\alpha_i = 0$  for all  $i$ , and  $E[MSB] = \sigma^2$  if and only if  $\beta_j = 0$  for all  $j$ .
- 11.4.** Consider a real-valued response variable  $Y$  and two independent variables  $X_1$  and  $X_2$ . Let  $L(Y|X_1)$ ,  $L(Y|X_2)$ , and  $L(Y|X_1, X_2)$  be the best linear predictors of  $Y$  given  $X_1$ ,  $Y$  given  $X_2$ , and  $Y$  given  $X_1, X_2$ , respectively. Partial correlation between  $Y$  and  $X_2$  given  $X_1$  is defined as  $\rho_{YX_2|X_1} = \text{Corr}[Y - L(Y|X_1), X_2 - L(X_2|X_1)]$ . Show that  $\rho_{YX_2|X_1}^2 = [E\{Y - L(Y|X_1)\}^2 - E\{Y - L(Y|X_1, X_2)\}^2]/E\{Y - L(Y|X_1)\}^2$ .
- 11.5.** Suppose a Gauss-Markov model is of the form  $\mathbf{Y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim N_n(0, \sigma^2 I)$ , where  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are of full rank, and the augmented matrix  $[\mathbf{X}_1 \mathbf{X}_2]$  is also of full rank. Let  $\mathbf{Q}_1 = \mathbf{X}_1^T(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T$  and  $\tilde{\mathbf{X}}_2 = (I - \mathbf{Q}_1)\mathbf{X}_2$ .
- (a) Show that the least squares estimate of  $\boldsymbol{\beta}_2$  is given by  $\hat{\boldsymbol{\beta}}_2 = (\tilde{\mathbf{X}}_2^T\tilde{\mathbf{X}}_2)^{-1}\tilde{\mathbf{X}}_2\mathbf{Y}$ .
  - (b) Find the distribution of  $\hat{\boldsymbol{\beta}}_2$  and use this to test  $H_0: \boldsymbol{\beta}_2 = 0$  vs  $H_1: \boldsymbol{\beta}_2 \neq 0$ .
- 11.6.** Consider a Gauss-Markov model  $\mathbf{Y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim N_n(0, \sigma^2 I)$ . We are interested testing  $H_0: \boldsymbol{\beta}_2 = 0$  vs  $H_1: \boldsymbol{\beta}_2 \neq 0$ . Let  $SSE_F$  and  $SSE_R$  be the residual sums of squares for the full and the reduced model (ie, the model under  $H_0$ ). The coefficient of partial determination is defined to be  $R_{YX_2|X_1}^2 = (SSE_R - SSE_F)/SSE_R$ .
- (a) Express  $R_{YX_2|X_1}^2$  as a function of the  $F$ -statistic for testing  $H_0: \boldsymbol{\beta}_2 = 0$  vs  $H_1: \boldsymbol{\beta}_2 \neq 0$ .
  - (b) Use the result in part (a) to describe the distribution of  $R_{YX_2|X_1}^2$  under  $H_0$ .  
[Hint: If  $U \sim \chi_p^2$ ,  $V \sim \chi_q^2$  and  $U$  and  $V$  are independent, then  $U/(U + V) \sim \text{Beta}(p/2, q/2)$ .]
- 11.7.** Assume that for a two-factor study with one observation for all treatment combinations, the appropriate model is  $Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$  with the usual constraints on  $\{\alpha_i\}$  and  $\{\beta_j\}$ . Here  $\{Y_{ij}\}$  are the observations and  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ .
- (a) Use the Scheffé method to obtain simultaneous confidence intervals for all contrasts in  $\{\alpha_i\}$ .

- (b) Use Tukey's method to obtain simultaneous confidence intervals for all pairwise difference of  $\{\alpha_i\}$ . How would you obtain simultaneous confidence intervals for all pairwise differences of  $\{\alpha_i\}$  and all pairwise differences of  $\{\beta_j\}$ , using Tukey's method, so that the family confidence is at least  $1 - \alpha$ ?

**11.8.** Prove [Lemma 11.7.2](#).

**11.9.** For the ANCOVA model with one factor and one covariate as given in

[Example 11.7.1](#), use the Scheffé method to obtain simultaneous confidence intervals for all contrasts of  $\{\alpha_i\}$ , where  $\{\alpha_i\}$  are the factor effects.

- 11.10.** Let  $\mathbf{Y}$  be  $n \times 1$  observation vector and assume that  $\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\varepsilon}$ , where the vector  $\boldsymbol{\varepsilon}$  consists of iid observations with mean 0 and variance  $\sigma^2$ . Consider a model of the form  $\mathbf{Y} = \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\varepsilon}$ , where  $\mathbf{X}_k$  is a  $n \times p_k$  matrix of rank  $p_k$ , which is being fitted to the data  $\mathbf{Y}$  and let  $\hat{\boldsymbol{\mu}}_k = \mathbf{X}_k \hat{\boldsymbol{\beta}}_k$  where  $\hat{\boldsymbol{\beta}}_k$  is the least squares estimator  $\boldsymbol{\beta}_k$ . The expected value of the squared distance between  $\boldsymbol{\mu}$  and  $\hat{\boldsymbol{\mu}}_k$  is  $D_k = E[\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_k\|^2]$ .

- (a) Let  $\hat{D}_k = FPE(k) - n\sigma^2$ , where the expression of Akaike's FPE is given in [Section 8.1](#) of this chapter. Is  $\hat{D}_k$  an unbiased estimate of  $D_k$ ? If not find its bias and find the condition under which  $\{E[\hat{D}_k] - D_k\}/D_k \rightarrow 0$  as  $n \rightarrow \infty$ .
- (b) Let  $\hat{D}_k = MAL(k) - n\sigma^2$ , where the expression of Mallows' criterion is given in [Section 8.1](#). Suppose we have a class of models  $\mathbf{Y} = \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\varepsilon}$ ,  $k = 1, \dots, K$ , where all the  $K - 1$  models are nested in the  $K$ th model (ie,  $\mathcal{M}(\mathbf{X}_k) \subset \mathcal{M}(\mathbf{X}_K)$ ). In Mallows' criterion let  $\hat{\sigma}^2$  be the MSE of the  $K$ th model. Is  $\hat{D}_k$  an unbiased estimate of  $D_k$ ? If not find its bias and find the condition under which  $\{E[\hat{D}_k] - D_k\}/D_k \rightarrow 0$  as  $n \rightarrow \infty$ .
- (c) In this part assume that  $\boldsymbol{\varepsilon} \sim N_n(0, \sigma^2 I)$ . Then the MLE of  $\sigma^2$  under the model  $\mathbf{Y} = \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\varepsilon}$  is  $\hat{\sigma}_k^2 = \|\mathbf{Y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k\|^2/n$ . If  $\tilde{\mathbf{Y}}$  is an independent copy of  $\mathbf{Y}$  but is independent of it, then the AIC is an estimate of  $-2E[\log f(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)]$ , where the expectation is over  $\mathbf{Y}$  and  $\tilde{\mathbf{Y}}$ , and  $\hat{\boldsymbol{\theta}}_k$  is the vector of  $\hat{\boldsymbol{\beta}}_k$  and  $\hat{\sigma}_k^2$  stacked vertically. Recall that  $AIC(k) = -2 \log f(\mathbf{Y}, \hat{\boldsymbol{\theta}}_k) + 2p_k$ , and denote  $-2 \log f(\tilde{\mathbf{Y}}, \hat{\boldsymbol{\theta}}_k)$  by  $L_k$ . If  $\boldsymbol{\mu} = \mathbf{X}_k \boldsymbol{\beta}_k$  for some  $\boldsymbol{\beta}_k$  (ie, the model being fitted is the correct one), then using asymptotic expansion as  $n \rightarrow \infty$  and  $p_k/n \rightarrow 0$ , obtain an approximation of  $E[AIC(k)] - E[L_k]$ , which is the bias of  $AIC(k)$  in estimating  $E[L_k]$ .

- 11.11.** Consider a model  $\mathbf{Y} = \sum_{j=1}^K \beta_j \mathbf{Z}_j + \boldsymbol{\varepsilon}$ , where  $\mathbf{Y}$  is  $n$ -dim observation vector, the vectors  $\{\mathbf{Z}_j\}$  are mutually orthogonal with  $\|\mathbf{Z}_j\|^2 = n$ , and  $\boldsymbol{\varepsilon} \sim N_n(0, \sigma^2 I)$ . Assume that  $\beta_j \neq 0$ ,  $j = 1, \dots, k^* < n$ , and  $\beta_j = 0$ ,  $j = k^* + 1, \dots, K$ . Consider the submodels  $\mathbf{Y} = \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\varepsilon}$ , where  $\mathbf{X}_k = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$  and  $\boldsymbol{\beta}_k = (\beta_1, \dots, \beta_k)^T$ ,  $k = 1, \dots, K$ . Consider a model selection criterion of the form  $F_k = \|\mathbf{Y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k\|^2 + ck\hat{\sigma}^2$ ,  $c > 0$ ,  $k = 1, \dots, K$ , where  $\hat{\boldsymbol{\beta}}_k$  is the least squares estimate of  $\boldsymbol{\beta}_k$  for the  $k$ th model and  $\hat{\sigma}^2$  is the MSE of the largest model under consideration (ie, the  $K$ th model). The values of  $F_1, F_2, \dots$  are calculated

sequentially and let  $\hat{k}$  be the index so that  $F_k - F_{k+1}$  becomes nonnegative for the first time (ie,  $F_k$  is strictly decreasing in  $k$ ,  $k = 1, \dots, \hat{k}$  and  $F_{\hat{k}+1} \geq F_{\hat{k}}$ ).

- (a) If we test  $H_{k0}: \beta_k = 0$  vs  $H_{k1}: \beta_k \neq 0$ ,  $k = 1, \dots, K$ , then consider the  $t$ -statistic  $t_k = \hat{\beta}_k / s(\hat{\beta}_k)$ , where  $\hat{\beta}_k = \mathbf{Z}_k^T \mathbf{Y}/n$  and  $s^2(\hat{\beta}_k) = \hat{\sigma}^2/n$ . Then given a critical value,  $H_{k0}$  is rejected or accepted depending on whether  $|t_k|$  is larger than the critical value or not. Show that  $|t_k| > c^{1/2}$  for  $1 \leq k \leq \hat{k}$  and  $|t_{\hat{k}+1}| \leq c^{1/2}$ , if model selection is done by using the criterion function  $\{F_k\}$ .
- (b) If  $c = \log n$ , then prove that  $P[\hat{k} = k^*] \rightarrow 1$  as  $n \rightarrow \infty$ . [This proves that a BIC-type criterion is capable of consistent model selection.]

**11.12.** Prove Lemma 11.10.1.

**11.13.** Prove Lemma 11.10.2.

**11.14.** In a one-factor random-effects model, compare the BLUP  $\tilde{\alpha}_1$  of  $\alpha_1$  to the naive predictor  $\hat{\alpha}_1 = \bar{Y}_{..} - \bar{Y}_{...}$ .

- (a) Obtain the distributions of  $\tilde{\alpha}_1 - \alpha_1$  and  $\hat{\alpha}_1 - \alpha_1$ .
- (b) Compare the errors  $E[(\tilde{\alpha}_1 - \alpha_1)^2]$  and  $E[(\hat{\alpha}_1 - \alpha_1)^2]$ .
- (c) Obtain the proportional reduction in the errors, that is,  $\theta = \left\{ E[(\hat{\alpha}_1 - \alpha_1)^2] - E[(\tilde{\alpha}_1 - \alpha_1)^2] \right\} / E[(\hat{\alpha}_1 - \alpha_1)^2]$  and examine it as  $\sigma_1^2 \rightarrow 0$  or  $\infty$  (assuming  $\sigma^2$  to be fixed).

**11.15.** In a repeated measures design each of the randomly selected  $m$  subjects (factor  $A$ ) is assigned to  $k$  levels of a treatment (factor  $B$ ). A reasonable model is thus

$$Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad j = 1, \dots, k, \quad i = 1, \dots, m,$$

where the subject effects  $\{\alpha_i\}$  are iid  $N(0, \sigma_\alpha^2)$ , treatment effects  $\{\beta_j\}$  are nonrandom with  $\sum \beta_j = 0$ ,  $\{\varepsilon_{ij}\}$  are iid  $N(0, \sigma^2)$ , and  $\{\alpha_i\}$  are independent of  $\{\varepsilon_{ij}\}$ . Define  $SSA = m \sum (\bar{Y}_{..} - \bar{Y}_{..})^2$ ,  $SSB = k \sum (\bar{Y}_{..j} - \bar{Y}_{..})^2$ , and  $SSE = \sum \sum (Y_{ij} - \bar{Y}_{..} - \bar{Y}_{..j} + \bar{Y}_{..})^2$ .

- (a) Find the mean, variance, and covariances of  $\{Y_{ij}\}$ .
- (b) Obtain  $E[SSA]$ ,  $E[SSB]$ , and  $E[SSE]$ .
- (c) Use the results in part (b) to obtain unbiased estimates of  $\sigma^2$  and  $\sigma_\alpha^2$ .
- (d) If someone ignores the subject effect and uses a model of the form  $Y_{ij} = \mu + \beta_j + \varepsilon_{ij}$  and obtains an MSE based on this model. Is this MSE an unbiased estimate of  $\sigma^2$ ? If not find its bias.

**11.16.** Consider the repeated measures design as in the previous exercise. It is of interest to obtain a BLUP of  $\alpha_1$ . The structure of the Mixed Model Equations suggest that the BLUP is a linear function of  $\{\bar{Y}_{..}\}$  and  $\{\bar{Y}_{..j}\}$ .

- (a) Assuming that  $\sigma_\alpha^2$  and  $\sigma^2$  are known, find the BLUP  $\tilde{\alpha}_1$  of  $\alpha_1$ .
- (b) Find the distribution of  $\tilde{\alpha}_1 - \alpha_1$ .

**11.17.** Verify the expressions for  $E[SSTR]$  and  $E[MSTR]$  given in Example 11.11.1.

**11.18.** Verify the expressions of  $E[MSA]$ ,  $E[MSB]$ , and  $E[MSAB]$  in Example 11.11.4.

- 11.19.** Verify the expressions of  $E[MSA]$ ,  $E[MSB]$ , and  $E[MSAB]$  in [Example 11.11.5](#).
- 11.20.** Consider a balanced two-factor ANOVA model in which factor  $A$  is fixed and factor  $B$  is random as in [Example 11.11.5](#). Let  $\mu_i$ ,  $\theta$ ,  $\hat{\mu}_i$ , and  $\hat{\theta}$  be as in that example.
- (a) Prove that  $(\hat{\mu}_i - \mu_i)/s(\hat{\mu}_i) \sim t_{a(b-1)}$ , where  $s^2(\hat{\mu}_i) = MSB(A)/(n_0 b)$ .
  - (b) Prove that  $(\hat{\theta} - \theta)/s(\hat{\theta}) \sim t_{(a-1)(b-1)}$ , where  $s^2(\hat{\theta}) = (\sum c_i^2)MSAB/(n_0 b)$ .

# Multivariate Analysis

## 12.1 Introduction

Multivariate analysis is an area of statistics which deals with observations that are vector valued. Almost all univariate statistical methods have their multivariate counterparts. For instance, when comparing two species of the same animal, various measures such as height, length, tail length, etc., may be measured. One can then compare these two species using a multivariate version of two-sample  $t$ -test. Fisher's famous Iris data set contains four measurements for each of the three species: petal length, petal width, sepal length, and sepal width. In order to compare the three species, a multivariate analog of analysis of variance has been developed. If there is a new observation vector (of unknown species) with four measurements, then allocation of this observation vector to one of the species is known as the problem of classification.

Another class of procedures has been developed for multivariate data which deal with dimensionality reduction. If many measurements are taken on children where each measurement is a measure of intelligence, then it is often the case that these various measures are correlated with each other. If there are 20 measurements for each child, it may be reasonable to look for a few summaries which contain most of the information. These summaries are often expressed as appropriate linear combinations of the measurements. This class of methods is known as principal components and factor analyses.

We describe these methods in a systematic manner starting with a few technical results on the Wishart distribution, which is a multivariate generalization of the chi-squared distribution.

## 12.2 Wishart Distribution

If  $\mathbf{Y}_1, \dots, \mathbf{Y}_k$  are iid  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$  and  $\mathbf{M} = \sum_{i=1}^k \mathbf{Y}_i \mathbf{Y}_i^T$ , then we say that  $\mathbf{M}$  has a (central) Wishart distribution with  $df = k$  and the scale matrix  $\boldsymbol{\Sigma}$ , and we write  $\mathbf{M} \sim W_p(k, \boldsymbol{\Sigma})$ . If the means  $\{\boldsymbol{\mu}_i\}$  of  $\{\mathbf{Y}_i\}$  are not necessarily equal to  $\mathbf{0}$ , then  $\mathbf{M}$  is said to have a non-central Wishart distribution  $W_p(k, \boldsymbol{\Sigma}, \boldsymbol{\Delta})$ , where  $\boldsymbol{\Delta} = (1/2) \boldsymbol{\Sigma}^{-1/2} \sum_{i=1}^k \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1/2}$  is the noncentrality matrix. Here we discuss only the central Wishart distribution and some of its properties. Some of the results stated below will be proved later in this chapter and further details on the theory can be found in the book by Mardia et al. [61]. In multivariate

analysis, we often deal with positive definite matrices and, for brevity of notation, we sometimes abbreviate “positive definite” by “pd.”

- (1) If  $\mathbf{M} \sim W_p(k, \Sigma)$ , then  $\mathbf{BMB}^T \sim W_m(k, \mathbf{B}\Sigma\mathbf{B}^T)$  for any  $m \times p$  matrix  $\mathbf{B}$ .
- (2) If  $\mathbf{M}_1 \sim W_p(k_1, \Sigma)$ ,  $\mathbf{M}_2 \sim W_p(k_2, \Sigma)$ , and  $\mathbf{M}_1$  and  $\mathbf{M}_2$  are independent, then  $\mathbf{M}_1 + \mathbf{M}_2 \sim W_p(k_1 + k_2, \Sigma)$ .
- (3) If  $\mathbf{M} \sim W_p(k, \Sigma)$  and  $\Sigma$  is pd, then  $P[\mathbf{M} \text{ is pd}] = 0$  if  $k < p$ , and  $P[\mathbf{M} \text{ is pd}] = 1$  if  $k \geq p$ .
- (4) If  $\mathbf{M} \sim W_p(k, \Sigma)$ , then  $\mathbf{a}^T \mathbf{M} \mathbf{a} / \mathbf{a}^T \Sigma \mathbf{a} \sim \chi_k^2$  if  $\mathbf{a}^T \Sigma \mathbf{a} \neq 0$ , where  $\mathbf{a}$  is in  $\mathbb{R}^p$ .
- (5) If  $\mathbf{M} \sim W_p(k, \Sigma)$ ,  $k \geq p$ , and  $\Sigma$  is pd, then  $\mathbf{a}^T \Sigma^{-1} \mathbf{a} / \mathbf{a}^T \mathbf{M}^{-1} \mathbf{a} \sim \chi_{k-p+1}^2$ , where  $\mathbf{a}$  is in  $\mathbb{R}^p$ .
- (6) Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $N_p(\mu, \Sigma)$ . Define

$$\bar{\mathbf{Y}} = \frac{1}{n} \sum \mathbf{Y}_i, \quad \mathbf{S} = \frac{1}{n-1} \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^T.$$

Then  $\bar{\mathbf{Y}}$  and  $\mathbf{S}$  are independent, and

$$\bar{\mathbf{Y}} \sim N_p(\mu, n^{-1} \Sigma), \quad (n-1)\mathbf{S} \sim W_p(n-1, \Sigma).$$

- (7) Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_{n_1}$  be iid  $N_p(\mu_1, \Sigma)$ ,  $\mathbf{Z}_1, \dots, \mathbf{Z}_{n_2}$  be iid  $N_p(\mu_2, \Sigma)$ , and assume that the samples  $\{\mathbf{Y}_i\}$  and  $\{\mathbf{Z}_j\}$  are independent. Define

$$\begin{aligned} \mathbf{S}_1 &= \frac{1}{n_1-1} \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^T, \\ \mathbf{S}_2 &= \frac{1}{n_2-1} \sum (\mathbf{Z}_j - \bar{\mathbf{Z}})(\mathbf{Z}_j - \bar{\mathbf{Z}})^T, \text{ and} \\ \mathbf{S}_{\text{pooled}} &= \frac{1}{n_1+n_2-2} [(n_1-1)\mathbf{S}_1 + (n_2-1)\mathbf{S}_2]. \end{aligned}$$

Then  $\bar{\mathbf{Y}} - \bar{\mathbf{Z}}$  and  $\mathbf{S}_{\text{pooled}}$  are independent, and

$$\begin{aligned} \bar{\mathbf{Y}} - \bar{\mathbf{Z}} &\sim N_p(\mu_1 - \mu_2, (1/n_1 + 1/n_2) \Sigma), \\ (n_1+n_2-2)\mathbf{S}_{\text{pooled}} &\sim W_p(n_1+n_2-2, \Sigma). \end{aligned}$$

- (8) Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n, \bar{\mathbf{Y}}$ , and  $\mathbf{S}$  be the same as in Property (6) above. The rv  $T^2 = n(\bar{\mathbf{Y}} - \mu)^T \mathbf{S}^{-1} (\bar{\mathbf{Y}} - \mu)$  is called Hotelling's  $T^2$ -statistic and it is distributed as  $\frac{(n-1)p}{n-p} F_{p,n-p}$ , where  $F_{p,n-p}$  has an  $F$ -distribution (central) with  $df = (p, n-p)$ .
- (9) Let  $\bar{\mathbf{Y}} - \bar{\mathbf{Z}}$  and  $\mathbf{S}_{\text{pooled}}$  be the same as in Property (7), and consider the following two-sample Hotelling's  $T^2$ -statistic

$$T^2 = (1/n_1 + 1/n_2)^{-1} (\bar{\mathbf{Y}} - \bar{\mathbf{Z}} - (\mu_1 - \mu_2))^T \mathbf{S}_{\text{pooled}}^{-1} (\bar{\mathbf{Y}} - \bar{\mathbf{Z}} - (\mu_1 - \mu_2)).$$

This two-sample Hotelling's  $T^2$  is distributed as  $\frac{(n_1+n_2-2)p}{n_1+n_2-p-1} F_{p,n_1+n_2-p-1}$ .

**(10)** If  $k \geq p$  and  $\Sigma$  is pd, the Wishart distribution has the pdf

$$f(\mathbf{M}) = \frac{|\mathbf{M}|^{(k-p-1)/2}}{2^{kp/2} \pi^{p(p-1)/4} |\Sigma|^{k/2}} \frac{\exp[-\text{trace}(\Sigma^{-1}\mathbf{M})/2]}{\prod_{i=1}^p \Gamma[(k+1-i)/2]}$$

where  $\mathbf{M}$  varies over pd matrices.

## 12.3 The Role of Multivariate Normal Distribution

We write  $\mathbf{Y} \sim (\mu, \Sigma)$  to mean that the  $p$ -dim random vector  $\mathbf{Y}$  has a mean  $\mu$  and covariance matrix  $\Sigma$ . Note that  $\mathbf{Y}$  is not necessarily normally distributed in this notation.

### 12.3.1 Mahalanobis Distance

If  $\mathbf{Y} \sim (\mu, \Sigma)$ , then the Mahalanobis distance between  $\mathbf{Y}$  and  $\mu$  is defined to be  $\Delta^2(\mathbf{Y}, \mu) = (\mathbf{Y} - \mu)^T \Sigma^{-1}(\mathbf{Y} - \mu)$ . Similarly, if  $\mathbf{Y}_1 \sim (\mu_1, \Sigma)$  and  $\mathbf{Y}_2 \sim (\mu_2, \Sigma)$ , then  $\Delta^2(\mathbf{Y}_1, \mathbf{Y}_2) = (\mathbf{Y}_1 - \mathbf{Y}_2)^T \Sigma^{-1}(\mathbf{Y}_1 - \mathbf{Y}_2)$ . Note that  $\Delta^2$  is well defined only if  $\Sigma$  is pd. It may be worthwhile to point out that the positive square root of  $\Delta^2$  is a distance on  $\mathbb{R}^p$  (and not  $\Delta^2$ ).

An important property of  $\Delta^2$  is that it is invariant under nonsingular linear transformations. Let  $\mathbf{X}_1 = \mathbf{a} + \mathbf{B}\mathbf{Y}_1$ ,  $\mathbf{X}_2 = \mathbf{a} + \mathbf{B}\mathbf{Y}_2$ , where  $\mathbf{a}$  is  $p \times 1$ ,  $\mathbf{B}$  is  $p \times p$  and is nonsingular. Then  $\Delta^2(\mathbf{X}_1, \mathbf{X}_2) = \Delta^2(\mathbf{Y}_1, \mathbf{Y}_2)$ . Mahalanobis distance comes up naturally in multivariate analysis. For instance, if  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  are iid  $(\mu, \Sigma)$ , then  $\Delta^2(\bar{\mathbf{Y}}, \mu) = n(\bar{\mathbf{Y}} - \mu)^T \Sigma^{-1}(\bar{\mathbf{Y}} - \mu)$ . If we want to test  $H_0: \mu = \mu_0$ , then we may use the Mahalanobis distance between  $\mathbf{Y}$  and  $\mu_0$ , that is,  $\Delta^2(\mathbf{Y}, \mu_0) = n(\bar{\mathbf{Y}} - \mu_0)^T \Sigma^{-1}(\bar{\mathbf{Y}} - \mu_0)$ , as a test statistic (assuming that  $\Sigma$  is known). If  $\Sigma$  is unknown and an estimate  $\hat{\Sigma}$  of  $\Sigma$  is available, then  $\Delta^2(\mathbf{Y}, \mu_0)$  can be approximated by  $n(\bar{\mathbf{Y}} - \mu_0)^T \hat{\Sigma}^{-1}(\bar{\mathbf{Y}} - \mu_0)$ .

In the univariate case we often assume normality. However, except for prediction intervals, almost all the inference are approximately valid without normality of the population as long as the sample size  $n$  is large. The same is also true in the multivariate case as long as  $n$  is large relative to  $p$ .

### 12.3.2 Multivariate Central Limit Theorem

If  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  are iid with mean  $\mu$  and covariance matrix  $\Sigma$ , then by the multivariate central limit theorem (Section A.4),

$$\sqrt{n}(\bar{\mathbf{Y}} - \mu) \xrightarrow{\mathcal{D}} N_p(\mathbf{0}, \Sigma) \text{ as } n \rightarrow \infty.$$

A consequence of this result is that

$$n(\bar{\mathbf{Y}} - \mu)^T \Sigma^{-1}(\bar{\mathbf{Y}} - \mu) \xrightarrow{\mathcal{D}} \chi_p^2 \text{ as } n \rightarrow \infty.$$

The sample covariance matrix  $\mathbf{S} = \frac{1}{n-1} \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^T$  is an unbiased estimate of  $\Sigma$ . Since  $\mathbf{S}$  is a consistent estimator of  $\Sigma$ , we have

$$T^2 = n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \mathbf{S}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\mu}) \xrightarrow{\mathcal{D}} \chi_p^2.$$

Property (8) in [Section 12.2](#) states that when the population is normal,

$$T^2 \sim \frac{(n-1)p}{n-p} F_{p,n-p}.$$

Since the rv  $F_{p,n-p}$  can be written as

$$F_{p,n-p} = \frac{W_1/p}{W_2/(n-p)},$$

where  $W_1 \sim \chi_p^2$ ,  $W_2 \sim \chi_{n-p}^2$ , and  $W_1$  and  $W_2$  are independent, we have

$$\frac{(n-1)p}{n-p} F_{p,n-p} = \frac{n-1}{n-p} \frac{W_1}{W_2/(n-p)}.$$

When  $n \rightarrow \infty$ ,  $(n-1)/(n-p) \rightarrow 1$  and  $W_2/(n-p) \xrightarrow{P} 1$ , and therefore

$$T^2 \sim \frac{(n-1)p}{n-p} F_{p,n-p} \xrightarrow{\mathcal{D}} \chi_p^2.$$

### 12.3.3 Checking Normality

A simple indication of multivariate normality is normality of each of the  $p$  component variables. Even though this may be enough in most cases, it is important to note that the normality of the marginal distributions does not imply multivariate normality. We now discuss a strategy for checking multivariate normality when we have iid  $p$ -dim observations  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  from a population.

If the population is indeed normal, then the quantities  $d_j^2 = (\mathbf{Y}_j - \bar{\mathbf{Y}})^T \mathbf{S}^{-1}(\mathbf{Y}_j - \bar{\mathbf{Y}})$  should be approximately iid  $\chi_p^2$ . So for a given data set, we can calculate these deviances  $\{d_j^2\}$  and plot them against the corresponding quantiles of the  $\chi_p^2$  distribution. If the population is multivariate normal, we expect the plot to be approximately linear. Here are the steps.

- (a) Order  $d_j^2$  from the smallest to the largest:  $d_{(1)}^2 \leq \dots \leq d_{(n)}^2$ .
- (b) Obtain the chi-squared plot, that is, plot  $\{d_{(j)}^2\}$  against  $\left\{ \chi_p^2((j-0.5)/n) \right\}$ , where  $\chi_p^2((j-0.5)/n)$  is the  $(j-0.5)/n$ -quantile of the  $\chi_p^2$  distribution.

We recommend the following steps for checking multivariate normality on the basis of a data set.

**Step I.** Check if each of the individual  $p$  variable is univariate normal.

**Step II.** Check if the chi-squared plot of the deviances  $\{d_j^2\}$  is linear.

Steps I and II do not guarantee multivariate normality. However, for practical purposes, these two steps are often enough for checking normality.

### 12.3.4 Sampling From a Normal Population

Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . As before, let  $\mathbf{S}$  (Property (6) in [Section 12.2](#)) be the sample covariance matrix. Then  $\bar{\mathbf{Y}}$  and  $\mathbf{S}$  are unbiased estimates of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , respectively. It turns out the maximum likelihood estimate (MLE) of  $\boldsymbol{\mu}$  is  $\bar{\mathbf{Y}}$ . However, the MLE of  $\boldsymbol{\Sigma}$  (proved below) is

$$\tilde{\mathbf{S}} = \frac{1}{n} \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^T = \frac{n-1}{n} \mathbf{S},$$

which is not an unbiased estimator for  $\boldsymbol{\Sigma}$ . We summarize the above and a bit more in the following result, the proof of which is given in [Section 12.3.5](#).

**Theorem 12.3.1.** *Let  $\bar{\mathbf{Y}}$  and  $\mathbf{S}$  be the sample mean and sample covariance matrix based on  $n$  independent observations from  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then:*

- (a) *The MLE of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are  $\bar{\mathbf{Y}}$  and  $\tilde{\mathbf{S}}$ , respectively.*
- (b) *Sufficient statistics for  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  are  $(\bar{\mathbf{Y}}, \mathbf{S})$ .*

### 12.3.5 Sampling Distributions

Results given in the following theorem are important for inference when sampling from a multivariate normal population.

**Theorem 12.3.2.** *Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid from  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then the following are true:*

- (a)  $\sqrt{n}(\bar{\mathbf{Y}} - \boldsymbol{\mu}) \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$ .
- (b)  $(n-1)\mathbf{S} \sim W_p(n-1, \boldsymbol{\Sigma})$ .
- (c)  $\bar{\mathbf{Y}}$  and  $\mathbf{S}$  are independent.
- (d)  $n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\mu}) \sim \chi_p^2$ .
- (e)  $T^2 = n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \mathbf{S}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\mu}) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$ .

*Proof of Theorem 12.3.1.*

- (a) The likelihood function  $L = L(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is given by

$$L = \left(1/\sqrt{2\pi}\right)^{np} (1/|\boldsymbol{\Sigma}|)^{n/2} \exp\left[-(1/2) \sum (\mathbf{Y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \boldsymbol{\mu})\right].$$

Maximizing  $L$  with respect to  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  is equivalent to minimizing  $-2 \log L$  with respect to  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Since

$$\begin{aligned} \sum (\mathbf{Y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}) &= \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}}) \\ &\quad + n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{Y}} - \boldsymbol{\mu}) \\ &= n \text{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) + n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{Y}} - \boldsymbol{\mu}), \end{aligned}$$

we have

$$-2 \log L = n \operatorname{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) + n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{Y}} - \boldsymbol{\mu}) + n \log(|\boldsymbol{\Sigma}|) + np \log(2\pi).$$

Clearly, if we minimize  $-2 \log L$  with respect to  $\boldsymbol{\mu}$ , the minimum occurs at  $\boldsymbol{\mu} = \bar{\mathbf{Y}}$ . So

$$-2 \log L(\bar{\mathbf{Y}}, \boldsymbol{\Sigma}) = n \operatorname{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) + n \log(|\boldsymbol{\Sigma}|) + np \log(2\pi).$$

In order to show that  $\tilde{\mathbf{S}}$  is indeed the MLE of  $\boldsymbol{\Sigma}$ , it is enough to show that the quantity  $\operatorname{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) + \log(|\boldsymbol{\Sigma}|)$  is minimized at  $\boldsymbol{\Sigma} = \tilde{\mathbf{S}}$ . Now

$$\begin{aligned} \operatorname{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) + \log(|\boldsymbol{\Sigma}|) &= \operatorname{trace}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) - \log(|\boldsymbol{\Sigma}^{-1}|) \\ &= \operatorname{trace}\left(\tilde{\mathbf{S}}^{1/2} \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}^{1/2}\right) - \log\left(\left|\tilde{\mathbf{S}}^{1/2} \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}^{1/2}\right|\right) + \log(|\tilde{\mathbf{S}}|) \\ &= \operatorname{trace}(\mathbf{R}) - \log(|\mathbf{R}|) + \log(|\tilde{\mathbf{S}}|), \end{aligned}$$

where  $\tilde{\mathbf{S}}^{1/2}$  is a symmetric square root of  $\tilde{\mathbf{S}}$  and  $\mathbf{R} = \tilde{\mathbf{S}}^{1/2} \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}^{1/2}$ . Let  $\lambda_1, \dots, \lambda_p$  be the eigenvalues of  $\mathbf{R}$ . Then

$$\operatorname{trace}(\mathbf{R}) - \log(|\mathbf{R}|) + \log(|\tilde{\mathbf{S}}|) = \sum \lambda_i - \sum \log \lambda_i + \log(|\tilde{\mathbf{S}}|).$$

The last quantity is minimized when  $\lambda_1 = \dots = \lambda_p = 1$ , that is,  $\mathbf{R} = \mathbf{I}$ . Consequently,  $\hat{\boldsymbol{\Sigma}} = \tilde{\mathbf{S}}$ .

**(b)** Follows from the factorization theorem. □

*Proof of Theorem 12.3.2.* The proof of part (a) is obvious. The proofs of parts (b) and (c) mirror their univariate counterparts. Part (d) follows from part (a). We now present the proof of the result in part (e).

Note that  $R_1 = n(\bar{\mathbf{Y}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{Y}} - \boldsymbol{\mu}) \sim \chi_p^2$ . Let  $\mathbf{d} = \sqrt{n}(\bar{\mathbf{Y}} - \boldsymbol{\mu})$ . Now we will use Property (5) of Wishart distribution from Section 12.2. Since  $\bar{\mathbf{Y}}$  and  $\mathbf{S}$  are independent, conditionally on  $\bar{\mathbf{Y}}$ ,

$$R_2 = \mathbf{d}^T \boldsymbol{\Sigma}^{-1} \mathbf{d} / [\mathbf{d}^T ((n-1)\mathbf{S})^{-1} \mathbf{d}] \sim \chi_{n-1-p+1}^2.$$

Since this conditional distribution does not depend on  $\bar{\mathbf{Y}}$ , we conclude that  $R_2 \sim \chi_{n-p}^2$  unconditionally, and,  $R_2$  is independent of  $\bar{\mathbf{Y}}$  and hence of  $R_1$ . Therefore,

$$T^2 = (n-1) \frac{R_1}{R_2} = \frac{(n-1)p}{n-p} \frac{R_1/p}{R_2/(n-p)} \sim \frac{(n-1)p}{n-p} F_{p,n-p}.$$

□

## 12.4 One-Sample Inference

Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . It is of interest to obtain the estimates of  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$ , and confidence regions for  $\boldsymbol{\mu}$ . In some cases we may be interested in testing

- I.  $H_0: \mu = \mu_0$  against  $H_1: \mu \neq \mu_0$  ( $\mu_0$  known).
- II.  $H_0: \psi = \psi_0$  against  $H_1: \psi \neq \psi_0$ , where  $\psi = D\mu$ ,  $D$  is an  $m \times p$  matrix of rank  $m \leq p$ , and  $D$  and  $\psi_0$  are known.

There are many examples of the second testing problem. A particular characteristic of the precision instruments produced by a company is considered to be important and the company takes a random sample of  $n$  instruments. The characteristic is measured by four engineers and thus there is a vector of four measurements for each instrument. If  $\mu = (\mu_1, \dots, \mu_4)^T$  is the mean vector, we may be interested in testing if these four measurements are same on the average, that is,  $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ . This can be restated as a testing problem given in (II) with

$$\psi = \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{pmatrix} = \begin{pmatrix} \mu_1 - \mu_4 \\ \mu_2 - \mu_4 \\ \mu_3 - \mu_4 \end{pmatrix} \text{ and } \psi_0 = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}.$$

### 12.4.1 Confidence Ellipsoid for $\mu$

Consider the random ellipsoid

$$A = \{\mathbf{u}: n(\bar{\mathbf{Y}} - \mathbf{u})^T S^{-1}(\bar{\mathbf{Y}} - \mathbf{u}) \leq c_\alpha\},$$

where  $c_\alpha = \frac{(n-1)p}{n-p} F_{p,n-p,\alpha}$  and  $F_{p,n-p,\alpha}$  is the  $(1 - \alpha)$ -quantile of the  $F$ -distribution with  $df = (p, n - p)$ . This method has been already discussed in the context of Linear models (Section 11.6.1 in Chapter 11). Since  $P_{\mu, \Sigma}[\mu \text{ is in } A] = 1 - \alpha$ , the set  $A$  is called the confidence ellipsoid for  $\mu$  with confidence coefficient  $1 - \alpha$ . When  $p = 2$  or  $3$ , it is not difficult to get a plot of this ellipsoid. However, it is not possible to visualize this ellipsoid when  $p \geq 4$ . In order to obtain confidence intervals for the individual components of  $\mu$ , one can consider the confidence shadows. However, this may sometimes lead to an inefficient method for simultaneous inference. We discuss two methods for constructing confidence intervals for linear combinations of  $\mu$ .

### 12.4.2 Simultaneous Confidence Intervals

Following the ideas given in Section 11.6.1 of Chapter 11, we present two methods for constructing simultaneous confidence intervals for linear combinations of the mean vector  $\mu$ . Proofs are not given since they are the same as in Chapter 11.

**(i) Scheffé method:** Simultaneous confidence intervals for all linear combinations  $\mathbf{l}^T \mu$ ,  $\mathbf{l}$  in  $\mathbb{R}^p$ , with a family confidence coefficient of  $1 - \alpha$  are

$$\mathbf{l}^T \mu: \mathbf{l}^T \bar{\mathbf{Y}} \pm \sqrt{c_\alpha} s(\mathbf{l}^T \bar{\mathbf{Y}}),$$

where  $s^2(\mathbf{l}^T \bar{\mathbf{Y}}) = \mathbf{l}^T S \mathbf{l} / n$  and  $c_\alpha = \frac{(n-1)p}{n-p} F_{p,n-p,\alpha}$ .

- (ii) *Bonferroni method:* Simultaneous confidence intervals for  $m$  linear combinations  $\mathbf{l}_1^T \boldsymbol{\mu}, \dots, \mathbf{l}_m^T \boldsymbol{\mu}$  with a family confidence coefficient at least  $1 - \alpha$  are

$$\mathbf{l}_i^T \boldsymbol{\mu} : \mathbf{l}_i^T \bar{\mathbf{Y}} \pm B s(\mathbf{l}_i^T \bar{\mathbf{Y}}), \quad i = 1, \dots, m,$$

where  $B = t_{n-1, \alpha/(2m)}$ .

As discussed in [Chapter 11](#), the Bonferroni method may be inefficient if  $m$  is not small since the multiplier  $B$  increases as  $m$  increases.

### 12.4.3 Hypothesis Testing

We now consider the two hypothesis testing problems involving the mean vector  $\boldsymbol{\mu}$  mentioned at the beginning of this section.

- I. Suppose we want to test  $H_0: \boldsymbol{\mu} = \boldsymbol{\mu}_0$  against  $H_1: \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$  ( $\boldsymbol{\mu}_0$  known). Then the test statistic is Hotelling's  $T^2$  statistic introduced earlier:  $T^2 = n(\bar{\mathbf{Y}} - \boldsymbol{\mu}_0)^T \mathbf{S}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\mu}_0)$ . Under  $H_0$ ,  $T^2 \sim \frac{(n-1)p}{n-p} F_{p, n-p}$ . We may reject  $H_0$  at level  $\alpha$  if  $T^2 > c_\alpha$ , where  $c_\alpha = \frac{(n-1)p}{n-p} F_{p, n-p, \alpha}$ .
- II. Let  $\boldsymbol{\psi} = \mathbf{D}\boldsymbol{\mu}$ , where  $\mathbf{D}$  is a known  $m \times p$  matrix of rank  $m \leq p$ . Suppose we want to test  $H_0: \boldsymbol{\psi} = \boldsymbol{\psi}_0$  against  $H_1: \boldsymbol{\psi} \neq \boldsymbol{\psi}_0$ , where  $\boldsymbol{\psi}_0$  is known. Then the appropriate test statistic is

$$T^2 = n(\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}_0)^T (\mathbf{D}\mathbf{S}\mathbf{D}^T)^{-1} (\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}_0).$$

Under  $H_0$ ,  $T^2 \sim \frac{(n-1)m}{n-m} F_{m, n-m}$ . So we reject  $H_0$  at level  $\alpha$  if  $T^2 > \frac{(n-1)m}{n-m} F_{m, n-m, \alpha}$ .

The second testing problem is the same as the first if we take the observation vectors to be  $\mathbf{W}_1 = \mathbf{D}\mathbf{Y}_1, \dots, \mathbf{W}_n = \mathbf{D}\mathbf{Y}_n$ . Note that  $\hat{\boldsymbol{\psi}} = \bar{\mathbf{W}} = \bar{\mathbf{D}\mathbf{Y}} \sim N_m(\mathbf{D}\boldsymbol{\mu}, \mathbf{D}\boldsymbol{\Sigma}\mathbf{D}^T)$  and an unbiased estimate of  $\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}^T$  is  $\mathbf{D}\mathbf{S}\mathbf{D}^T$ , where  $\mathbf{S}$  is the sample covariance matrix based on the sample  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ .

### 12.4.4 Likelihood Ratio Test

Let us denote  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  by  $\boldsymbol{\theta}$ . The likelihood function is

$$L(\boldsymbol{\theta}) = \left(1/\sqrt{2\pi}\right)^{np} (1/|\boldsymbol{\Sigma}|)^{n/2} \exp\left[-(1/2) \sum (\mathbf{Y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \boldsymbol{\mu})\right].$$

The likelihood ratio statistic for testing  $H_0$  is

$$\lambda = \frac{\max_{\boldsymbol{\theta} \text{ in } H_0} L(\boldsymbol{\theta})}{\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta})} = \frac{\max\{L(\boldsymbol{\theta}): \boldsymbol{\theta} \text{ in the reduced model}\}}{\max\{L(\boldsymbol{\theta}): \boldsymbol{\theta} \text{ in the full model}\}},$$

where the “reduced model” is the same as the “model under  $H_0$ .” Asymptotic theory for likelihood ratio tests tells us, under  $H_0$ ,

$$-2 \log \lambda \xrightarrow{\mathcal{D}} \chi_t^2$$

as  $n \rightarrow \infty$ , where

$$\begin{aligned} t &= (\# \text{ of parameters estimated under the full model}) \\ &\quad - (\# \text{ of parameters estimated under } H_0) \\ &= p. \end{aligned}$$

For large  $n$ , we reject  $H_0$  at level  $\alpha$  if  $-2 \log \lambda > \chi^2_{p,\alpha}$ . It should be noted that the chi-square approximation for  $-2 \log \lambda$  is valid asymptotically, and it is possible to obtain the exact distribution of  $-2 \log \lambda$  using [Theorem 12.4.1](#) stated below.

What is the relation between the likelihood ratio test statistic and Hotelling's  $T^2$ ? The following result provides the answer when the true mean  $\mu$  is in a small neighborhood of  $\mu_0$ .

**Theorem 12.4.1.** Consider the problem of testing  $H_0: \mu = \mu_0$  vs  $H_1: \mu \neq \mu_0$  on the basis of an iid sample  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  from  $N_p(\mu, \Sigma)$ .

- (a) The likelihood ratio test statistic for this testing problem is  $\lambda = |\tilde{\mathbf{S}}|^{n/2} / |\mathbf{S}_0|^{n/2}$ , where  $\mathbf{S}_0 = n^{-1} \sum (\mathbf{Y}_i - \mu_0)(\mathbf{Y}_i - \mu_0)^T$  and  $\tilde{\mathbf{S}}$  is as given in [Theorem 12.3.1](#).
- (b) Let  $T^2 = n(\bar{\mathbf{Y}} - \mu_0)^T \mathbf{S}^{-1} (\bar{\mathbf{Y}} - \mu_0)$ . Under  $P_{\mu, \Sigma}$ , where  $\mu$  is in the set

$$A_n = \{\mathbf{u}: \|\mathbf{u} - \mu_0\| \leq c_n\}, c_n \rightarrow 0 \text{ as } n \rightarrow \infty,$$

$$-2 \log \lambda = T^2 + o_P(1).$$

*Remark 12.4.1.* An examination of the proof of [Theorem 12.4.1](#) shows that if  $c_n$  is of order  $n^{-1/2}$ , then the  $o_P(1)$  in the last theorem can be replaced by  $O_P(1/n)$ . Hence, under  $H_0: \mu = \mu_0$ , or under a contiguous alternative of the form  $\mu = \mu_0 + n^{-1/2}\delta$ , we have  $-2 \log \lambda = T^2 + O_P(1/n)$ .

The following important identity will be used in the proof of [Theorem 12.4.1](#)

$$\mathbf{S}_0 = \tilde{\mathbf{S}} + (\bar{\mathbf{Y}} - \mu_0)(\bar{\mathbf{Y}} - \mu_0)^T,$$

where  $\mathbf{S}_0$  is as given in the theorem.

*Proof of Theorem 12.4.1.*

- (a) We use the expression for the likelihood function  $L(\mu, \Sigma)$  given in the proof of [Theorem 12.3.1](#). Note that the MLEs of  $\mu$  and  $\Sigma$  are obtained by maximizing  $L(\mu, \Sigma)$  over  $\mu$  and  $\Sigma$ . Since the MLEs of  $\mu$  and  $\Sigma$  are  $\bar{\mathbf{Y}}$  and  $\tilde{\mathbf{S}}$ , respectively, we have

$$\begin{aligned} \max_{\mu, \Sigma} L(\mu, \Sigma) &= L(\bar{\mathbf{Y}}, \tilde{\mathbf{S}}) \\ &= \left(1/\sqrt{2\pi}\right)^{np} \left(1/|\tilde{\mathbf{S}}|\right)^{n/2} \exp\left[-(1/2) \sum (\mathbf{Y}_i - \bar{\mathbf{Y}})^T \tilde{\mathbf{S}}^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}})\right] \\ &= \left(1/\sqrt{2\pi}\right)^{np} \left(1/|\tilde{\mathbf{S}}|\right)^{n/2} \exp(-pn/2). \end{aligned}$$

When  $\mu = \mu_0$ , we may write the likelihood as

$$\begin{aligned} L(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) &= \left(1/\sqrt{2\pi}\right)^{np} (1/|\boldsymbol{\Sigma}|)^{n/2} \exp\left[-(1/2) \sum (\mathbf{Y}_i - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}_0)\right] \\ &= \left(1/\sqrt{2\pi}\right)^{npnp} (1/|\boldsymbol{\Sigma}|)^{n/2} \exp\left[-(n/2) \text{trace}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right]. \end{aligned}$$

The arguments (Proof of [Theorem 12.3.1](#)) employed in obtaining the MLE of  $\boldsymbol{\Sigma}$ , when  $\boldsymbol{\mu}$  is unknown, can be used when  $\boldsymbol{\mu} = \boldsymbol{\mu}_0$ . Basically the same arguments show that  $L(\boldsymbol{\mu}_0, \boldsymbol{\Sigma})$  is maximized at  $\boldsymbol{\Sigma} = \mathbf{S}_0$  where

$$\mathbf{S}_0 = \frac{1}{n} \sum (\mathbf{Y}_i - \boldsymbol{\mu}_0)(\mathbf{Y}_i - \boldsymbol{\mu}_0)^T.$$

Consequently,

$$\begin{aligned} \max_{\boldsymbol{\mu}=\boldsymbol{\mu}_0, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \max_{\boldsymbol{\Sigma}} L(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) \\ &= \left(1/\sqrt{2\pi}\right)^{np} (1/|\mathbf{S}_0|)^{n/2} \exp\left(-(n/2) \text{trace}(\mathbf{S}_0^{-1} \mathbf{S}_0)\right) \\ &= \left(1/\sqrt{2\pi}\right)^{np} (1/|\mathbf{S}_0|)^{n/2} \exp(-np/2). \end{aligned}$$

Hence the likelihood ratio statistic for the testing problem is

$$\lambda = \frac{\max_{\boldsymbol{\mu}=\boldsymbol{\mu}_0, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}, \boldsymbol{\Sigma})} = \frac{|\tilde{\mathbf{S}}|^{n/2}}{|\mathbf{S}_0|^{n/2}}.$$

- (b)** Using the identity stated before the beginning of the proof of this theorem, we can express the likelihood ratio statistic  $\lambda$  as

$$\lambda = \frac{1}{|\tilde{\mathbf{S}}^{-1/2} \mathbf{S}_0 \tilde{\mathbf{S}}^{-1/2}|} = \frac{1}{|\mathbf{I} + \mathbf{b} \mathbf{b}^T|},$$

where  $\mathbf{b} = \tilde{\mathbf{S}}^{-1/2}(\bar{\mathbf{Y}} - \boldsymbol{\mu}_0)$ .

It is not difficult to check that the matrix  $\mathbf{I} + \mathbf{b} \mathbf{b}^T$  has an eigenvalue equal to 1 with multiplicity  $p - 1$ , and the remaining eigenvalue is  $1 + \|\mathbf{b}\|^2$ . Since the determinant of a matrix is the product of its eigenvalues, we get

$$|\mathbf{I} + \mathbf{b} \mathbf{b}^T| = 1 + \|\mathbf{b}\|^2 = 1 + (\bar{\mathbf{Y}} - \boldsymbol{\mu}_0)^T \tilde{\mathbf{S}}^{-1} (\bar{\mathbf{Y}} - \boldsymbol{\mu}_0) = 1 + \frac{1}{n-1} T^2.$$

Since  $\mathbf{S} = \boldsymbol{\Sigma} + o_P(1)$  and  $\bar{\mathbf{Y}} = \boldsymbol{\mu} + O_P(n^{-1/2})$  (as  $n \rightarrow \infty$ ),

$$T^2 \leq n \|\bar{\mathbf{Y}} - \boldsymbol{\mu}_0\|^2 \|\mathbf{S}\| \leq n \left[ \|(\bar{\mathbf{Y}} - \boldsymbol{\mu})\| + \|(\boldsymbol{\mu} - \boldsymbol{\mu}_0)\| \right]^2 \|\mathbf{S}\| = o_P(n)$$

on the set  $A_n = \{\boldsymbol{\mu}: \|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| \leq c_n\}$ , and this implies  $T^2/(n-1) = o_P(1)$  on the set  $A_n$ . Hence

$$-2 \log \lambda = n \log(1 + T^2/(n-1)) = T^2 + o_P(1).$$

□

## 12.5 Two-Sample Problem

Suppose that we have two independent samples from two multivariate normal populations with different mean vectors, but the same covariance matrix. Let  $\mathbf{Y}_{1j}, j = 1, \dots, n_1$ , be iid  $N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$  and  $\mathbf{Y}_{2j}, j = 1, \dots, n_2$  be iid  $N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$ . We assume that the samples  $\{\mathbf{Y}_{1j}\}$  and  $\{\mathbf{Y}_{2j}\}$  are independent. We address the following two issues:

- (a) test for  $H_0: \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$  against  $H_1: \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ .
- (b) confidence statements for  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ .

### 12.5.1 Estimation

The MLEs for  $\boldsymbol{\mu}_1$  and  $\boldsymbol{\mu}_2$  are

$$\hat{\boldsymbol{\mu}}_1 = \bar{\mathbf{Y}}_{1\cdot} = (1/n_1) \sum \mathbf{Y}_{1j}, \text{ and } \hat{\boldsymbol{\mu}}_2 = \bar{\mathbf{Y}}_{2\cdot} = (1/n_2) \sum \mathbf{Y}_{2j}.$$

An unbiased estimate of  $\boldsymbol{\Sigma}$  is

$$\mathbf{S} = \mathbf{S}_{\text{pooled}} = (n_1 + n_2 - 2)^{-1} [(n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2],$$

where  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are the sample covariance matrices on the basis of  $\{\mathbf{Y}_{1j}: j = 1, \dots, n_1\}$  and  $\{\mathbf{Y}_{2j}: j = 1, \dots, n_2\}$ , respectively, that is,

$$\mathbf{S}_i = (n_i - 1)^{-1} \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i\cdot})(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i\cdot})^T, \quad i = 1, 2.$$

As in the univariate two-sample problem, the  $\mathbf{S} = \mathbf{S}_{\text{pooled}}$  is a better estimator of  $\boldsymbol{\Sigma}$  than  $\mathbf{S}_1$  or  $\mathbf{S}_2$ .

The following result is useful for inference in two-sample problems.

**Theorem 12.5.1.** *Let us denote  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$  and its MLE  $\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2$  by  $\boldsymbol{\delta}$  and  $\hat{\boldsymbol{\delta}}$ , respectively. The following hold:*

- (a)  $(\hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\mu}}_2, \mathbf{S})$  is sufficient for  $(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$ .
- (b)  $\hat{\boldsymbol{\mu}}_1 \sim N_p(\boldsymbol{\mu}_1, (1/n_1)\boldsymbol{\Sigma})$ ,  $\hat{\boldsymbol{\mu}}_2 \sim N_p(\boldsymbol{\mu}_2, (1/n_2)\boldsymbol{\Sigma})$ , and  $\hat{\boldsymbol{\delta}} \sim N_p(\boldsymbol{\delta}, (1/n_1 + 1/n_2)\boldsymbol{\Sigma})$ .
- (c)  $(n_1 + n_2 - 2)\mathbf{S} \sim W_p(n_1 + n_2 - 2, \boldsymbol{\Sigma})$ .
- (d)  $\hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\mu}}_2$ , and  $\mathbf{S}$  are independent.
- (e) The two-sample Hotelling  $T^2$ -statistic

$$T^2 = (1/n_1 + 1/n_2)^{-1} (\hat{\boldsymbol{\delta}} - \boldsymbol{\delta})^T \mathbf{S}^{-1} (\hat{\boldsymbol{\delta}} - \boldsymbol{\delta})$$

has the same distribution as  $\frac{(n_1+n_2-2)p}{n_1+n_2-p-1} F_{p, n_1+n_2-p-1}$ .

When  $n_1 + n_2 \rightarrow \infty$ , the two-sample Hotelling's  $T^2 \xrightarrow{\mathcal{D}} \chi_p^2$ .

*Proof of Theorem 12.5.1.* We only provide a proof of parts (c) and (e) since the rest are not difficult to establish.

- (c) Note that  $(n_1 - 1)\mathbf{S}_1 \sim W_p(n_1 - 1, \boldsymbol{\Sigma})$ ,  $(n_2 - 1)\mathbf{S}_2 \sim W_p(n_2 - 1, \boldsymbol{\Sigma})$  and that  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are independent. Hence

$$(n_1 + n_2 - 2)\mathbf{S} = (n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2 \sim W_p(n_1 + n_2 - 2, \boldsymbol{\Sigma}).$$

- (e) We will argue as in the one-sample case. Let

$$\mathbf{d} = (1/n_1 + 1/n_2)^{-1/2}(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}),$$

$$R_1 = (1/n_1 + 1/n_2)^{-1}(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta})^T \boldsymbol{\Sigma}^{-1}(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}), \text{ and}$$

$$R_2 = \mathbf{d}^T \boldsymbol{\Sigma}^{-1} \mathbf{d} / \{\mathbf{d}^T ((n_1 + n_2 - 2)\mathbf{S})^{-1} \mathbf{d}\}.$$

From part (a),  $R_1 \sim \chi_p^2$ . From part (d),  $\hat{\boldsymbol{\mu}}_1$  and  $\hat{\boldsymbol{\mu}}_2$  are independent of  $\mathbf{S}$  and hence conditionally on  $\hat{\boldsymbol{\mu}}_1$  and  $\hat{\boldsymbol{\mu}}_2$ ,  $R_2 \sim \chi_{n_1+n_2-p-1}^2$  (by Property (5) in [Section 12.2](#)). Since this conditional distribution does not depend on  $\hat{\boldsymbol{\mu}}_1$  and  $\hat{\boldsymbol{\mu}}_2$ , we conclude that  $R_2 \sim \chi_{n_1+n_2-p-1}^2$  unconditionally, and  $R_1$  and  $R_2$  are independent. Therefore,

$$\begin{aligned} T^2 &= (n_1 + n_2 - 2) \frac{R_1}{R_2} \\ &= \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} \frac{R_1/p}{R_2/(n_1 + n_2 - p - 1)} \\ &\sim \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}. \end{aligned}$$

□

### 12.5.2 Hypothesis Testing

We want to test  $H_0: \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$  against  $H_1: \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$  at a level of significance  $\alpha$ . Consider the following Hotelling's  $T^2$ -statistic

$$T^2 = (1/n_1 + 1/n_2)^{-1}(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2)^T \mathbf{S}^{-1}(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2).$$

Under  $H_0$ ,  $T^2$  has the same distribution as  $\frac{(n_1+n_2-2)p}{n_1+n_2-p-1} F_{p, n_1 + n_2 - p - 1}$ . So we reject  $H_0$  at level  $\alpha$  if  $T^2 > c_\alpha$  where

$$c_\alpha = \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1, \alpha}.$$

### 12.5.3 Confidence Ellipsoid for $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$

Consider the ellipsoid

$$A = \{\boldsymbol{\delta}: (1/n_1 + 1/n_2)^{-1}(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2 - \boldsymbol{\delta})^T \mathbf{S}^{-1}(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2 - \boldsymbol{\delta}) \leq c_\alpha\}.$$

Since  $P_{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}}[\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 \text{ is in } A] = 1 - \alpha$ ,  $A$  is a confidence ellipsoid for  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$  with confidence coefficient  $1 - \alpha$ .

### 12.5.4 Simultaneous Confidence Intervals

- (a) *Scheffé method*: Simultaneous confidence intervals for all linear combinations  $\mathbf{l}^T \boldsymbol{\delta}$ ,  $\mathbf{l}$  in  $\mathbb{R}^p$  and  $\boldsymbol{\delta} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ , with a family confidence coefficient  $1 - \alpha$  are given by  $\mathbf{l}^T \hat{\boldsymbol{\delta}} \pm \sqrt{c_\alpha} s(\mathbf{l}^T \hat{\boldsymbol{\delta}})$ , where  $\hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2$  and  $s^2(\mathbf{l}^T \hat{\boldsymbol{\delta}}) = (1/n_1 + 1/n_2) \mathbf{l}^T \mathbf{S} \mathbf{l}$ .
- (b) *Bonferroni method*: Simultaneous confidence intervals for  $\mathbf{l}_1^T \boldsymbol{\delta}, \dots, \mathbf{l}_m^T \boldsymbol{\delta}$  with a family confidence coefficient of at least  $1 - \alpha$  are

$$\mathbf{l}_i^T \boldsymbol{\delta}: \mathbf{l}_i^T \hat{\boldsymbol{\delta}} \pm t_{n_1+n_2-2,\alpha/(2m)} s(\mathbf{l}_i^T \hat{\boldsymbol{\delta}}), \quad i = 1, \dots, m.$$

## 12.6 One-Factor MANOVA

Suppose that we have  $k$  multivariate normal populations with possibly different mean vectors, but the same covariance matrix. Let  $\{\mathbf{Y}_{ij}: j = 1, \dots, n_i\}$  be iid  $N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ ,  $i = 1, \dots, k$ . We may write the one-factor MANOVA model as

$$\mathbf{Y}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\varepsilon}_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\{\boldsymbol{\varepsilon}_{ij}\}$  are iid  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$ . We can also rewrite the above as a factor-effect model

$$\mathbf{Y}_{ij} = \boldsymbol{\mu} + \boldsymbol{\alpha}_i + \boldsymbol{\varepsilon}_{ij},$$

where  $\boldsymbol{\mu} = \sum(n_i/n)\boldsymbol{\mu}_i$ ,  $\boldsymbol{\alpha}_i = \boldsymbol{\mu}_i - \boldsymbol{\mu}$  and  $n = n_1 + \dots + n_k$  is the total number of observation vectors.

The following issues are of interest:

- (a) test  $H_0: \boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_k$  against  $H_1$ : not all  $\boldsymbol{\mu}_i$ 's are the same.  
 (b) confidence statements about  $\boldsymbol{\mu}_i$ 's and  $\boldsymbol{\alpha}_i$ 's.

### 12.6.1 Estimation

MLEs for  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$  are  $\hat{\boldsymbol{\mu}}_i = \bar{\mathbf{Y}}_{i..} = (1/n_i) \sum_j \mathbf{Y}_{ij}$ ,  $i = 1, \dots, k$ . An unbiased estimate of the covariance matrix  $\boldsymbol{\Sigma}$  is

$$\begin{aligned} \mathbf{S} &= \frac{1}{n-k} [(n_1-1)\mathbf{S}_1 + \dots + (n_k-1)\mathbf{S}_k], \text{ where} \\ \mathbf{S}_i &= \frac{1}{n_i-1} \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i..})(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{i..})^T \end{aligned}$$

is the sample covariance matrix on the basis of the  $i$ th sample,  $i = 1, \dots, k$ . The MLEs for  $\boldsymbol{\mu}$  and  $\boldsymbol{\alpha}_i$  are

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{Y}}_{...}, \quad \hat{\boldsymbol{\alpha}}_i = \hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\mu}} = \bar{\mathbf{Y}}_{i..} - \bar{\mathbf{Y}}_{...}$$

As in one- and two-sample cases, the MLE of  $\boldsymbol{\Sigma}$

$$\tilde{\mathbf{S}} = \frac{1}{n} [(n_1-1)\mathbf{S}_1 + \dots + (n_k-1)\mathbf{S}_k]$$

is a biased estimator.

The following result is useful for inference in one-factor MANOVA.

**Theorem 12.6.1.** *Let  $\hat{\mu}_1, \dots, \hat{\mu}_k$  be the sample means based on independent random samples from  $N_p(\mu_1, \Sigma), \dots, N_p(\mu_k, \Sigma)$ , and let  $S$  be the pooled covariance matrix. Then the following hold:*

- (a)  $(\hat{\mu}_1, \dots, \hat{\mu}_k, S)$  is sufficient for  $(\mu_1, \dots, \mu_k, \Sigma)$ .
- (b)  $\bar{Y}_i \sim N_p(\mu_i, (1/n_i)\Sigma)$ .
- (c)  $\hat{\mu} \sim N_p(\mu, (1/n)\Sigma)$ .
- (d)  $\hat{\alpha}_i \sim N_p(\alpha_i, (1/n_i - 1/n)\Sigma)$ .
- (e)  $\hat{\mu}_1, \dots, \hat{\mu}_k$  and  $S$  are all independent.
- (f)  $(n - k)S \sim W_p(n - k, \Sigma)$ .
- (g) When  $\mu_1 = \dots = \mu_k$ ,  $\sum n_i (\bar{Y}_i - \bar{Y}_{..})(\bar{Y}_i - \bar{Y}_{..})^T \sim W_p(k - 1, \Sigma)$ .

The proof of this theorem is not given. The results in it can be obtained using arguments similar to one- and two-sample cases given above, and by borrowing appropriate results from univariate analysis of variance.

## 12.6.2 Hypothesis Testing in One-Factor MANOVA

Suppose we wish to test  $H_0: \mu_1 = \dots = \mu_k$  against  $H_1: \text{not all } \mu_i\text{'s are the same}$ . This is equivalent to testing  $H_0: \alpha_1 = \dots = \alpha_k = \mathbf{0}$  against  $H_1: \text{not all } \alpha_i\text{'s are equal to } \mathbf{0}$ . We define a few matrices analogous to the various sums of squares in the univariate case. Total sum of squares and products (SSP), between group SSP and within group SSP, as well as their corresponding degrees of freedom are given below

$$\begin{aligned}\mathbf{T} &= \sum_i \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..})(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..})^T, \quad df = n - 1, \\ \mathbf{B} &= \sum_i n_i (\bar{\mathbf{Y}}_i - \bar{\mathbf{Y}}_{..})(\bar{\mathbf{Y}}_i - \bar{\mathbf{Y}}_{..})^T, \quad df = k - 1, \text{ and} \\ \mathbf{W} &= \sum_i \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)^T, \quad df = n - k.\end{aligned}$$

It is fairly easy to check that  $\mathbf{T} = \mathbf{B} + \mathbf{W}$ .

**Theorem 12.6.2.** *Consider the problem of testing  $H_0: \alpha_1 = \dots = \alpha_k = \mathbf{0}$  against  $H_1: \text{not all } \alpha_i\text{'s are equal to } \mathbf{0}$ , in one-factor MANOVA. The likelihood ratio test statistic is*

$$\lambda = \left\{ \frac{|\mathbf{W}|}{|\mathbf{T}|} \right\}^{n/2} = \Lambda^{n/2},$$

where  $\Lambda = \frac{|\mathbf{W}|}{|\mathbf{T}|}$  is called Wilks' lambda.

The proof of this result will be given below in [Section 12.6.4](#), but we describe the inference procedures first. The exact distribution of the likelihood ratio statistic  $\lambda$  is complicated even under  $H_0$ . In some cases we know the exact distributional results as given below. However, for the general case, we may use the asymptotic theory. From the theory of likelihood ratio tests, under  $H_0$ ,

$$-2 \log \lambda = -n \log \Lambda \xrightarrow{\mathcal{D}} \chi^2_{p(k-1)}$$

as  $n \rightarrow \infty$ . In order to obtain a better asymptotic approximation, Bartlett made a suitable modification to the test statistic. According to his modification, under  $H_0$ ,

$$-[n - 1 - (p + k)/2] \log \Lambda \xrightarrow{\mathcal{D}} \chi^2_{p(k-1)},$$

as  $n \rightarrow \infty$ . If  $\alpha$  is the given level of significance, we may reject  $H_0$  at level  $\alpha$  if

$$-[n - 1 - (p + k)/2] \log \Lambda > \chi^2_{p(k-1), \alpha}.$$

### 12.6.3 Simultaneous Confidence Intervals

If we are interested in constructing simultaneous confidence intervals for  $\alpha_{il}$ ,  $l = 1, \dots, p$ ,  $i = 1, \dots, k$ , then the Bonferroni method yields the intervals

$$\alpha_{il}: \hat{\alpha}_{il} \pm t_{n-k, \alpha/(2pk)} s(\hat{\alpha}_{il}),$$

where  $\hat{\alpha}_{il}$  is the  $l$ th component of  $\hat{\alpha}_i$  and  $s^2(\hat{\alpha}_{il}) = (1/n_i - 1/n)s_{ll}$  and  $s_{ll}$  is the  $l$ th diagonal element of  $\mathbf{S}$ . Similarly if we want to construct simultaneous confidence intervals for all pairwise differences  $\alpha_{il} - \alpha_{i'l}$ ,  $1 \leq i \neq i' \leq k$ ,  $l = 1, \dots, p$ , with a family confidence coefficient of at least  $1 - \alpha$ , then the intervals are

$$\alpha_{il} - \alpha_{i'l}: \hat{\alpha}_{il} - \hat{\alpha}_{i'l} \pm t_{n-k, \alpha/[pk(k-1)]} s(\hat{\alpha}_{il} - \hat{\alpha}_{i'l}),$$

where  $\hat{\alpha}_{il} - \hat{\alpha}_{i'l}$  is the  $l$ th component of  $\hat{\alpha}_i - \hat{\alpha}_{i'} = \bar{\mathbf{Y}}_i - \bar{\mathbf{Y}}_{i'}$ , and

$$s^2(\hat{\alpha}_{il} - \hat{\alpha}_{i'l}) = (1/n_i + 1/n_{i'})s_{ll}.$$

Clearly, the Bonferroni method may lead to wide intervals if  $p$  (and/or  $k$ ) is not small. If  $k$  is not large, we may use one-factor ANOVA models to construct Scheffé- or Tukey-type (in the balanced case) confidence intervals for  $\alpha_{il} - \alpha_{i'l}$ ,  $1 \leq i \neq i' \leq k$ , for each  $l = 1, \dots, p$ , so that the family confidence coefficient is at least  $1 - \alpha$ .

### 12.6.4 Exact Distributions of Wilks' Lambda

In general it is not easy to obtain the exact distribution of Wilks'  $\Lambda$  statistic and one usually uses the asymptotic distribution with Bartlett corrections. However, there are some special cases where exact distributions have been obtained and some of these are presented in following table:

$p = 1$	$k \geq 2$	$\frac{n-k}{k-1} \frac{1-\Lambda}{\Lambda} \sim F_{k-1, n-k}$
$p = 2$	$k \geq 2$	$\frac{n-k-1}{k-1} \frac{1-\sqrt{\Lambda}}{\sqrt{\Lambda}} \sim F_{2(k-1), 2(n-k-1)}$
$p \geq 1$	$k = 2$	$\frac{n-p-1}{p} \frac{1-\Lambda}{\Lambda} \sim F_{p, n-p-1}$
$p \geq 1$	$k = 3$	$\frac{n-p-2}{p} \frac{1-\sqrt{\Lambda}}{\sqrt{\Lambda}} \sim F_{2p, 2(n-p-2)}$

*Proof of Theorem 12.6.2.* Note that the likelihood is

$$L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) = \left(1/\sqrt{2\pi}\right)^{pn} (1/|\boldsymbol{\Sigma}|)^{n/2} \exp \left[ - \sum_i \sum_j (\mathbf{Y}_{ij} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_{ij} - \boldsymbol{\mu}_i)/2 \right].$$

Under  $H_0: \boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_k$ , the MLE for the common mean  $\boldsymbol{\mu}$  is  $\hat{\boldsymbol{\mu}} = \bar{\mathbf{Y}}_{..}$  and the MLE for  $\boldsymbol{\Sigma}$  is

$$\mathbf{S}_0 = (1/n) \sum_i \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..})(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_{..})^T.$$

For the general case, the MLE for  $\boldsymbol{\mu}_i$  is  $\hat{\boldsymbol{\mu}}_i = \bar{\mathbf{Y}}_i$  and the MLE of  $\boldsymbol{\Sigma}$  is

$$\tilde{\mathbf{S}} = (1/n) \sum_i \sum_j (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)(\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)^T.$$

Following the argument similar to the one- and two-sample cases, we can show that

$$\begin{aligned} \max_{\boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_k, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) &= L(\bar{\mathbf{Y}}_{..}, \mathbf{S}_0) \\ &= \left(1/\sqrt{2\pi}\right)^{pn} (1/|\mathbf{S}_0|)^{n/2} \exp(-pn/2), \text{ and} \\ \max_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) &= L(\bar{\mathbf{Y}}_1, \dots, \bar{\mathbf{Y}}_k, \tilde{\mathbf{S}}) \\ &= \left(1/\sqrt{2\pi}\right)^{pn} (1/|\tilde{\mathbf{S}}|)^{n/2} \exp(-pn/2). \end{aligned}$$

So the likelihood ratio test statistic for testing equality of the means is

$$\lambda = \frac{\max_{\boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_k, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma})}{\max_{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma})} = \left( \frac{|\tilde{\mathbf{S}}|}{|\mathbf{S}_0|} \right)^{n/2} = \left( \frac{|\mathbf{W}|}{|\mathbf{T}|} \right)^{n/2}.$$

The last equality holds since  $\mathbf{S}_0 = (1/n)\mathbf{T}$  and  $\tilde{\mathbf{S}} = (1/n)\mathbf{W}$ .

□

### 12.6.5 More Tests for One-Factor MANOVA

In Section 12.6.2, we have discussed the likelihood ratio test for  $H_0: \boldsymbol{\alpha}_1 = \dots = \boldsymbol{\alpha}_k = \mathbf{0}$  against  $H_1$ : not all  $\boldsymbol{\alpha}_i$ 's are  $\mathbf{0}$ . There are three other well-known tests, and computer packages routinely report them. These are

- (a) Lawley-Hotelling trace:  $\text{trace}(\mathbf{B}\mathbf{W}^{-1})$ ,
- (b) Pillai's trace:  $\text{trace}(\mathbf{B}\mathbf{T}^{-1})$ , and
- (c) Roy's largest root: the largest eigenvalue of  $\mathbf{B}\mathbf{T}^{-1}$ .

For each of these tests, we reject  $H_0$  if the corresponding statistic is larger than a threshold value. It turns out that all these four test statistics can be written as functions of the eigenvalues of  $\mathbf{B}\mathbf{T}^{-1}$ . Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$  be the eigenvalues of  $\mathbf{B}$  with respect to  $\mathbf{T}$  (ie, the eigenvalues of  $\mathbf{T}^{-1/2}\mathbf{B}\mathbf{T}^{-1/2}$ , where  $\mathbf{T}^{-1/2}$  is symmetric). Then we have

- (a) likelihood ratio:  $\left\{ \prod_{i=1}^p (1 - \hat{\lambda}_i) \right\}^{n/2}$ ,
- (b) Lawley-Hotelling trace:  $\sum_{i=1}^p \hat{\lambda}_i / (1 - \hat{\lambda}_i)$ ,
- (c) Pillai's trace:  $\sum_{i=1}^p \hat{\lambda}_i$ , and
- (d) Roy's largest root:  $\hat{\lambda}_1$ .

*Remark 12.6.1.* Since the rank of the matrix  $\mathbf{B}$  is  $s = \min(p, k-1)$ , the number of nonzero generalized eigenvalues of  $\mathbf{B}$  with respect to  $\mathbf{T}$  is equal to  $s$ . Thus  $\hat{\lambda}_j > 0$ ,  $j = 1, \dots, s$ , and the remaining  $\hat{\lambda}_j$ 's are 0.

### Interpretation of Tests in MANOVA

When testing  $H_0: \boldsymbol{\alpha}_1 = \dots = \boldsymbol{\alpha}_k = \mathbf{0}$ , the reduced and the full MANOVA models are

$$\begin{aligned} Y_{ij} &= \mu + \boldsymbol{\epsilon}_{ij} \text{ (reduced),} \\ Y_{ij} &= \mu + \boldsymbol{\alpha}_i + \boldsymbol{\epsilon}_{ij} \text{ (full).} \end{aligned}$$

Now if we look at  $\{\mathbf{e}^T Y_{ij}\}$ ,  $\mathbf{e} \in \mathbb{R}^p$ , then we get the reduced and full models for the one-factor ANOVA case

$$\begin{aligned} \mathbf{e}^T Y_{ij} &= \mathbf{e}^T \mu + \mathbf{e}^T \boldsymbol{\epsilon}_{ij} \text{ (reduced),} \\ \mathbf{e}^T Y_{ij} &= \mathbf{e}^T \mu + \mathbf{e}^T \boldsymbol{\alpha}_i + \mathbf{e}^T \boldsymbol{\epsilon}_{ij} \text{ (full).} \end{aligned}$$

Since the coefficient of determination  $R^2(\mathbf{e})$  is the proportional reduction in the residual sum of squares from the reduced to the full model (following the terminology used in Chapter 11), we have

$$R^2(\mathbf{e}) = \frac{SSE_R - SSE_F}{SSE_R} = \frac{\sum n_i (\mathbf{e}^T \bar{Y}_{..} - \mathbf{e}^T \bar{Y}_{..})^2}{\sum_i \sum_j (\mathbf{e}^T Y_{ij} - \mathbf{e}^T \bar{Y}_{..})^2} = \frac{\mathbf{e}^T \mathbf{B} \mathbf{e}}{\mathbf{e}^T \mathbf{T} \mathbf{e}}. \quad (1)$$

Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$  be the generalized eigenvalues of  $\mathbf{B}$  with respect to  $\mathbf{T}$ , that is,  $\hat{\lambda}_1, \dots, \hat{\lambda}_p$  are the eigenvalues of  $\mathbf{T}^{-1/2} \mathbf{B} \mathbf{T}^{-1/2}$  with the corresponding eigenvectors  $\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_p$  (Section B.4). If we denote  $\hat{\mathbf{e}}_i = \mathbf{T}^{-1/2} \hat{\mathbf{u}}_i$ , then  $R^2(\hat{\mathbf{e}}_i) = \hat{\lambda}_i$ ,  $i = 1, \dots, p$ .

Clearly,  $R^2(\mathbf{e})$  is maximized at  $\mathbf{e} = \hat{\mathbf{e}}_1$  and  $R^2(\hat{\mathbf{e}}_1) = \hat{\lambda}_1$ . The next largest value of  $R^2(\mathbf{e})$  is obtained by maximizing it over  $\mathbf{e}$  subject to the constraint  $\mathbf{e}^T \mathbf{T} \hat{\mathbf{e}}_1 = 0$ . This maximum is attained at  $\mathbf{e} = \hat{\mathbf{e}}_2$  and  $R^2(\hat{\mathbf{e}}_2) = \lambda_2$ . This argument can be carried out further and it shows that  $\{R^2(\hat{\mathbf{e}}_i)\}$  are the same as the generalized eigenvalues  $\{\hat{\lambda}_i\}$  of  $\mathbf{B}$  with respect to  $\mathbf{T}$ .

We can now express the four test statistics given above for the one-factor MANOVA in terms of  $R^2(\hat{\mathbf{e}}_1), \dots, R^2(\hat{\mathbf{e}}_p)$ ,

- (a) likelihood ratio:  $\left\{ \prod_{i=1}^p (1 - R^2(\hat{\mathbf{e}}_i)) \right\}^{n/2}$ ,
- (b) Lawley-Hotelling trace:  $\sum_{i=1}^p R^2(\hat{\mathbf{e}}_i) / [1 - R^2(\hat{\mathbf{e}}_i)]$ ,
- (c) Pillai trace:  $\sum_{i=1}^p R^2(\hat{\mathbf{e}}_i)$ , and
- (d) Roy's largest root:  $R^2(\hat{\mathbf{e}}_1)$ .

## 12.7 Two-Factor MANOVA

Let us consider a two-factor balanced MANOVA model

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk},$$

$k = 1, \dots, n_0, j = 1, \dots, b, i = 1, \dots, a$ , where  $\{\epsilon_{ijk}\}$  are iid  $N_p(\mathbf{0}, \Sigma)$ ,

$$\begin{aligned} \sum \alpha_i &= \mathbf{0}, \quad \sum \beta_j = \mathbf{0}, \\ \sum_i (\alpha\beta)_{ij} &= \mathbf{0} \text{ for all } j, \text{ and } \sum_j (\alpha\beta)_{ij} = \mathbf{0} \text{ for all } i. \end{aligned}$$

For notational simplicity we write  $\gamma_{ij}$  for  $(\alpha\beta)_{ij}$ . The vectors  $\{\alpha_i\}$ ,  $\{\beta_j\}$ , and  $\{\gamma_{ij}\}$  are the main effects of factor A, main effects of factor B, and interaction effects, respectively. The total number of observation vectors is  $n = n_0ab$ .

### 12.7.1 Estimation

The MLEs of  $\mu$ ,  $\alpha_i$ ,  $\beta_j$ , and  $\gamma_{ij}$  are

$$\begin{aligned} \hat{\mu} &= \bar{Y}_{...}, \quad \hat{\alpha}_i = \bar{Y}_{i..} - \bar{Y}_{...}, \quad \hat{\beta}_j = \bar{Y}_{.j} - \bar{Y}_{...}, \text{ and} \\ \hat{\gamma}_{ij} &= \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j} + \bar{Y}_{...}. \end{aligned}$$

We now write down the sums of squares and products matrices along with their degrees of freedom,

$$\begin{aligned} \mathbf{SSP}_{\text{tot}} &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{...})(Y_{ijk} - \bar{Y}_{...})^T, \quad df = n - 1 \\ \mathbf{SSP}_A &= n_0 b \sum_i \hat{\alpha}_i \hat{\alpha}_i^T, \quad df = a - 1, \\ \mathbf{SSP}_B &= n_0 a \sum_j \hat{\beta}_j \hat{\beta}_j^T, \quad df = b - 1, \\ \mathbf{SSP}_{AB} &= n_0 \sum_i \sum_j \hat{\gamma}_{ij} \hat{\gamma}_{ij}^T, \quad df = (a - 1)(b - 1), \\ \mathbf{SSP}_{\text{res}} &= \sum_i \sum_j \sum_k (Y_{ijk} - \bar{Y}_{ij.})(Y_{ijk} - \bar{Y}_{ij.})^T, \quad df = n - ab. \end{aligned}$$

The following lemma lists some useful facts.

**Lemma 12.7.1.** *For a two-factor MANOVA, the following hold:*

- (a)  $\mathbf{SSP}_{\text{tot}} = \mathbf{SSP}_A + \mathbf{SSP}_B + \mathbf{SSP}_{AB} + \mathbf{SSP}_{\text{res}}$ .
- (b)  $df(\mathbf{SSP}_{\text{tot}}) = df(\mathbf{SSP}_A) + df(\mathbf{SSP}_B) + df(\mathbf{SSP}_{AB}) + df(\mathbf{SSP}_{\text{res}})$ .
- (c)  $\mathbf{SSP}_A$ ,  $\mathbf{SSP}_B$ ,  $\mathbf{SSP}_{AB}$ , and  $\mathbf{SSP}_{\text{res}}$  are independent.
- (d)  $\mathbf{SSP}_{\text{res}} \sim W_p(n - ab, \Sigma)$ .
- (e) An unbiased estimate of  $\Sigma$  is  $\frac{1}{n-ab} \mathbf{SSP}_{\text{res}}$ .

### 12.7.2 Hypothesis Testing in Two-Factor MANOVA

We will write down here only the likelihood ratio tests. However, as in the one-factor MANOVA case, there are other tests such as those by Pillai, Lawley-Hotelling, and Roy, and computer packages routinely report them.

#### *Test for Interactions*

Suppose we wish to test  $H_0: \boldsymbol{\gamma}_{ij} = \mathbf{0}$  for all  $i$  and  $j$ , against  $H_1$ : not all  $\boldsymbol{\gamma}_{ij}$  are zero.

Wilks' lambda for this test is  $\Lambda = \frac{|\text{SSP}_{\text{res}}|}{|\text{SSP}_{AB} + \text{SSP}_{\text{res}}|}$  and the likelihood ratio test statistic is  $\lambda = \Lambda^{n/2}$ . Under  $H_0$ , using Bartlett's modification we have

$$-\left\{ab(n_0 - 1) - \frac{p + 1 - (a - 1)(b - 1)}{2}\right\} \log \Lambda \xrightarrow{\mathcal{D}} \chi^2_{(a-1)(b-1)p},$$

as  $n \rightarrow \infty$ . We can use this result to obtain the critical value or the  $p$ -value in order to carry out the test.

#### *Test for the Main Effects of Factor A*

If we wish to test  $H_0: \boldsymbol{\alpha}_i = \mathbf{0}$  for all  $i$ , against  $H_1$ : not all  $\boldsymbol{\alpha}_i$  are zero, Wilks' lambda criterion is  $\Lambda = \frac{|\text{SSP}_{\text{res}}|}{|\text{SSP}_A + \text{SSP}_{\text{res}}|}$  and the likelihood ratio test statistic is  $\lambda = \Lambda^{n/2}$ .

Once again, we can use the following asymptotic result to carry out this test. Under  $H_0$ , using Bartlett's approximation we have

$$-\left\{ab(n_0 - 1) - \frac{p + 1 - (a - 1)}{2}\right\} \log \Lambda \xrightarrow{\mathcal{D}} \chi^2_{(a-1)p},$$

as  $n \rightarrow \infty$ .

#### *Test for the Main Effects of Factor B*

If we wish to test  $H_0: \boldsymbol{\beta}_j = \mathbf{0}$  for all  $j$ , against  $H_1$ : not all  $\boldsymbol{\beta}_j$  are zero, Wilks' lambda criterion is  $\Lambda = \frac{|\text{SSP}_{\text{res}}|}{|\text{SSP}_B + \text{SSP}_{\text{res}}|}$  and the likelihood ratio test statistic is  $\lambda = \Lambda^{n/2}$ . It turns out that under  $H_0$  (using Bartlett's modification),

$$-\left\{ab(n_0 - 1) - \frac{p + 1 - (b - 1)}{2}\right\} \log \Lambda \xrightarrow{\mathcal{D}} \chi^2_{(b-1)p},$$

as  $n \rightarrow \infty$ .

### 12.7.3 Simultaneous Confidence Intervals

Simultaneous confidence intervals with a family confidence coefficient of at least  $1 - \alpha$  for all the pairwise differences in the main effect of factor  $A$  are given by

$$\alpha_{il} - \alpha_{i'l}: \hat{\alpha}_{il} - \hat{\alpha}_{i'l} + Bs(\hat{\alpha}_{il} - \hat{\alpha}_{i'l}), \quad 1 \leq i \neq i' \leq a, \quad l = 1, \dots, p,$$

where,  $B = t_{n-ab,\alpha/[pa(a-1)]}$ ,  $s(\hat{\alpha}_{il} - \hat{\alpha}_{i'l}) = \sqrt{2s_{ll}/(n_0 b)}$ , and  $s_{ll}$  is the  $l$ th diagonal element of  $\mathbf{S}$ .

Simultaneous confidence intervals with a family confidence coefficient of at least  $1 - \alpha$  for all the pairwise differences in the main effect of factor  $B$  are given by

$$\beta_{jl} - \beta_{j'l}: \hat{\beta}_{jl} - \hat{\beta}_{j'l} + Bs(\hat{\beta}_{jl} - \hat{\beta}_{j'l}), \quad 1 \leq j \neq j' \leq b, \quad l = 1, \dots, p,$$

where  $B = t_{n-ab,\alpha/[pb(b-1)]}$  and  $s(\hat{\beta}_{jl} - \hat{\beta}_{j'l}) = \sqrt{2s_{ll}/(n_0a)}$ .

Simultaneous confidence intervals with a family confidence coefficient of at least  $1 - \alpha$  for all the pairwise differences in the mean response are given by

$$\begin{aligned} \mu_{ijl} - \mu_{i'j'l}: & \hat{\mu}_{ijl} - \hat{\mu}_{i'j'l} + Bs(\hat{\mu}_{ijl} - \hat{\mu}_{i'j'l}), \\ & 1 \leq i \neq i' \leq a, \quad 1 \leq j \neq j' \leq b, \quad (i, j) \neq (i', j'), \quad l = 1, \dots, p, \end{aligned}$$

where  $B = t_{n-ab,\alpha/[pab(ab-1)]}$  and  $s(\hat{\mu}_{ijl} - \hat{\mu}_{i'j'l}) = \sqrt{2s_{ll}/n_0}$ . [Note that  $\hat{\mu}_{ij} = \bar{Y}_{ij..}$ .]

The simultaneous confidence intervals given above are expected to be wide. Therefore, we may use Tukey's method for pairwise comparisons for each of the  $p$  univariate ANOVA models with a confidence level of  $1 - \alpha/p$  so that the overall confidence level is at least  $1 - \alpha$ .

## 12.8 Multivariate Linear Model

In the usual linear model framework, the response for each of the  $n$  observations is real valued. Thus the observed vector of responses  $\mathbf{Y}$  is  $n$ -dim, and if the design matrix  $\mathbf{X}$  is  $n \times k$ , then the Gauss-Markov model is written as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\beta}$  is an unknown vector of parameters and  $\boldsymbol{\epsilon}$  is an  $n \times 1$  vector of mutually uncorrelated random errors with mean zero and common variance  $\sigma^2$  ([Section 11.1 in Chapter 11](#)). If each of the  $n$  observations is  $p$ -dim, then we have an  $n \times p$  matrix  $\mathbf{Y}$  of observed responses which can modeled by a generalization of the framework described in [Section 11.1 of Chapter 11](#) resulting in a multivariate linear model. We have also seen in [Chapter 11](#) that the Gauss-Markov setup and its extensions include regression, analysis of variance, analysis of covariance, random- and mixed-effect models as special cases, and the same is true for their multivariate counterparts.

In the multivariate case, we have  $p$  columns of  $n \times 1$  observation vectors  $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ , and for each observed vector of responses, the model is

$$\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \boldsymbol{\epsilon}_j, \quad j = 1, \dots, p,$$

where  $\boldsymbol{\beta}_j$  is a  $k$ -dim unknown vector of parameters and  $\boldsymbol{\epsilon}_j$  is an  $n$ -dim vector of mutually uncorrelated mean zero random errors. Thus we have

$$\begin{aligned} [\mathbf{Y}_1, \dots, \mathbf{Y}_p] &= \mathbf{X}[\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_p] + [\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_p], \text{ or} \\ \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \end{aligned}$$

where  $\mathbf{Y}$  is  $n \times p$  with columns  $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ ,  $\mathbf{X}$  is  $n \times k$ ,  $\boldsymbol{\beta}$  is  $k \times p$  with columns  $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_p$ , and  $\boldsymbol{\epsilon}$  is  $n \times p$  with columns  $\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_p$ . This is the multivariate linear model and it has the same formal structure as in Eq. (1) in [Chapter 11](#). It is possible to analyze each of the  $p$  linear models separately, but this strategy has a drawback since the  $p$  observations in each

row of the response matrix may be correlated, that is, the  $p$  elements in each row of  $\boldsymbol{\epsilon}$  may be correlated, and a procedure which analyzes the  $p$  linear models separately, fails to take into account this dependence. A joint analysis of these  $p$  models is therefore preferable. In the subsequent discussion we assume that  $\mathbf{X}^T \mathbf{X}$  is nonsingular.

In the multivariate linear model, we often assume that the rows of  $\boldsymbol{\epsilon}$  are iid  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$ . The goal is to estimate the matrices of unknown parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$ , carry out inferential procedures such as tests of hypotheses and construction of confidence intervals, and make predictions whenever necessary. We do not discuss random- and mixed-effect cases here. We write down a few basic results on the estimation of  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$ , and the distributions of the estimates.

### 12.8.1 Estimation of $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$

The normal equations and least squares estimate of  $\boldsymbol{\beta}$  are similar to those in the univariate case, and they are

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{Y}, \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

In order to see why the same equations come up, we may consider rewriting the model. Postmultiplying both sides of the multivariate linear model by  $\boldsymbol{\Sigma}^{-1/2}$  leads to

$$\begin{aligned} [\mathbf{R}_1, \dots, \mathbf{R}_p] &= \mathbf{X} [\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p] + [\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_p], \text{ or} \\ \mathbf{R} &= \mathbf{X} \boldsymbol{\gamma} + \boldsymbol{\delta}, \end{aligned}$$

where  $\mathbf{R}_j = \mathbf{Y}_j \boldsymbol{\Sigma}^{-1/2}$ ,  $\boldsymbol{\gamma}_j = \boldsymbol{\beta}_j \boldsymbol{\Sigma}^{-1/2}$ ,  $\boldsymbol{\delta}_j = \boldsymbol{\epsilon}_j \boldsymbol{\Sigma}^{-1/2}$ ,  $j = 1, \dots, p$ . Under normality (ie, the rows of  $\boldsymbol{\epsilon}$  are iid  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$ ), the error vectors  $\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_p$  are mutually independent, and each  $\boldsymbol{\delta}_j$  has  $n$  entries which are iid with mean 0 and variance 1. Even if the assumption of normality is not valid,  $\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_p$  are mutually uncorrelated, and we can minimize

$$\sum \| \mathbf{R}_j - \mathbf{X} \boldsymbol{\gamma}_j \|^2$$

with respect to  $\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p$  in order to get the least squares estimates, which lead to the following  $p$  normal equations

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \boldsymbol{\gamma}_j &= \mathbf{X}^T \mathbf{R}_j, \quad \text{ie, } \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}_j \boldsymbol{\Sigma}^{-1/2} = \mathbf{X}^T \mathbf{Y}_j \boldsymbol{\Sigma}^{-1/2}, \quad \text{ie,} \\ \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}_j &= \mathbf{X}^T \mathbf{Y}_j, \quad j = 1, \dots, p, \text{ or} \\ \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} &= \mathbf{X}^T \mathbf{Y}, \quad \text{ie, } \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \end{aligned}$$

Since

$$\begin{aligned} \sum \| \mathbf{R}_j - \mathbf{X} \boldsymbol{\gamma}_j \|^2 &= \text{trace}((\mathbf{R} - \mathbf{X} \boldsymbol{\gamma})^T (\mathbf{R} - \mathbf{X} \boldsymbol{\gamma})) \\ &= \text{trace}((\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) \boldsymbol{\Sigma}^{-1}), \end{aligned}$$

it follows that the normal equations are obtained by minimizing  $\text{trace}((\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) \boldsymbol{\Sigma}^{-1})$  with respect to  $\boldsymbol{\beta}$ .

An unbiased estimate of  $\Sigma$  is given by

$$\mathbf{S} = \frac{1}{n-k}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

As in the univariate case, the estimated response  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  and the residuals  $\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  are uncorrelated, and are independent when the rows of  $\boldsymbol{\epsilon}$  are iid  $N_p(\mathbf{0}, \Sigma)$ , in which case,  $\mathbf{S}$  and  $\hat{\boldsymbol{\beta}}$  are independent as well.

### 12.8.2 Properties of the Estimates of $\boldsymbol{\beta}$ and $\Sigma$

For the rest of this section, we assume that the rows of  $\boldsymbol{\epsilon}$  are iid  $N_p(\mathbf{0}, \Sigma)$ . We have the following result which can be used for tests of hypotheses and construction of confidence intervals.

**Theorem 12.8.1.**

(a) The MLEs of  $\boldsymbol{\beta}$  and  $\Sigma$  are

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}, \quad \tilde{\mathbf{S}} = \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

- (b)  $(\hat{\boldsymbol{\beta}}, \mathbf{S})$  are sufficient for  $(\boldsymbol{\beta}, \Sigma)$ .
- (c)  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  are independent.
- (d)  $\hat{\boldsymbol{\beta}}$  and  $\mathbf{S}$  are independent.
- (e) For any vector  $\mathbf{b}$  in  $\mathbb{R}^p$ ,  $\hat{\boldsymbol{\beta}}\mathbf{b} \sim N_k(\boldsymbol{\beta}\mathbf{b}, \mathbf{b}^T \Sigma \mathbf{b} (\mathbf{X}^T \mathbf{X})^{-1})$ .
- (f)  $(n-k)\mathbf{S} \sim W_p(n-k, \Sigma)$ .

## 12.9 Principal Components Analysis

Principal components analysis is a widely used method in multivariate analysis, and its goal is to reduce the dimensionality of the data with as little loss of information as possible. We first describe the basic ideas behind principal components, and discuss estimation issues later. Let  $\mathbf{Y}$  be a  $p$ -dim random vector with mean  $\mu$  and covariance matrix  $\Sigma$ . In some cases we assume that the diagonal elements of  $\Sigma$  are 1 which means that  $\Sigma$  is a correlation matrix.

If  $\{\mathbf{e}_1, \dots, \mathbf{e}_p\}$  is an orthonormal basis of  $\mathbb{R}^p$ , then

$$\begin{aligned} \mathbf{Y} - \mu &= \sum_{j=1}^p [\mathbf{e}_j^T (\mathbf{Y} - \mu)] \mathbf{e}_j, \text{ and} \\ \|\mathbf{Y} - \mu\|^2 &= \sum_{j=1}^p [\mathbf{e}_j^T (\mathbf{Y} - \mu)]^2. \end{aligned}$$

Note that  $\mathbf{e}_j^T (\mathbf{Y} - \mu)$  is a random variable with mean 0 and variance  $\mathbf{e}_j^T \Sigma \mathbf{e}_j$ . The total variability of  $\mathbf{Y}$  is

$$\begin{aligned}\sum_{j=1}^p \text{Var}[Y_j] &= \sum_{j=1}^p \text{E}[Y_j - \mu_j]^2 = \text{E}\left\{\sum_{j=1}^p (Y_j - \mu_j)^2\right\} \\ &= \text{E}[\|\mathbf{Y} - \boldsymbol{\mu}\|^2] = \text{trace}(\boldsymbol{\Sigma}).\end{aligned}$$

On the other hand

$$\sum_{j=1}^p \text{Var}[Y_j] = \text{E}[\|\mathbf{Y} - \boldsymbol{\mu}\|^2] = \sum_{j=1}^p \text{E}[\mathbf{e}_j^T (\mathbf{Y} - \boldsymbol{\mu})]^2 = \sum_{j=1}^p \text{Var}[\mathbf{e}_j^T \mathbf{Y}].$$

Now if it happens that for some  $k < p$  (and hopefully  $k$  is small),

$$\text{E}[\|\mathbf{Y} - \boldsymbol{\mu}\|^2] \approx \sum_{j=1}^k \text{E}[\mathbf{e}_j^T (\mathbf{Y} - \boldsymbol{\mu})]^2, \text{ and}$$

$$\sum_{j=k+1}^p \text{E}[\mathbf{e}_j^T (\mathbf{Y} - \boldsymbol{\mu})]^2 \text{ is small,}$$

then the total variability of  $\mathbf{Y}$  is explainable (to a large extent) by the variability of  $k$  random variables  $\mathbf{e}_1^T \mathbf{Y}, \dots, \mathbf{e}_k^T \mathbf{Y}$ . Ideally, the reduction in dimensionality is substantial if  $k$  is small in comparison to  $p$ . Let  $\lambda_1 \geq \dots \geq \lambda_p$  be the eigenvalues of  $\boldsymbol{\Sigma}$  with  $\mathbf{u}_1, \dots, \mathbf{u}_p$ , the corresponding orthonormal eigenvectors. From the properties of eigenvalues and eigenvectors (Section B.2) we have that

$$\max\left\{\mathbf{e}_1^T \boldsymbol{\Sigma} \mathbf{e}_1 + \dots + \mathbf{e}_k^T \boldsymbol{\Sigma} \mathbf{e}_k : \mathbf{e}_1, \dots, \mathbf{e}_k \text{ orthonormal}\right\} = \lambda_1 + \dots + \lambda_k,$$

and the maximum is attained at  $\mathbf{e}_1 = \mathbf{u}_1, \dots, \mathbf{e}_k = \mathbf{u}_k$ . So

$$\begin{aligned}\max\left\{\text{Var}[\mathbf{e}_1^T \mathbf{Y}] + \dots + \text{Var}[\mathbf{e}_k^T \mathbf{Y}] : \mathbf{e}_1, \dots, \mathbf{e}_k \text{ orthonormal}\right\} \\ = \text{Var}[\mathbf{u}_1^T \mathbf{Y}] + \dots + \text{Var}[\mathbf{u}_k^T \mathbf{Y}] = \lambda_1 + \dots + \lambda_k.\end{aligned}$$

*Remark 12.9.1.*

- (a) The random variable  $\mathbf{u}_1^T (\mathbf{Y} - \boldsymbol{\mu})$  is called the first principal component of  $\mathbf{Y}$ ,  $\mathbf{u}_2^T (\mathbf{Y} - \boldsymbol{\mu})$  is the second principal component, and so on. The vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots$  are sometimes called the loading vectors.
- (b) Principal components and the eigenvalues are not unit free. For this reason, in many cases one may prefer to carry out principal components analysis on the standardized variables which means that  $\boldsymbol{\Sigma}$  is the correlation matrix.
- (c) Note that  $\pi_k = (\lambda_1 + \dots + \lambda_k)/(\lambda_1 + \dots + \lambda_p)$  is the proportion of the total variability of  $\mathbf{Y}$  explained by the first  $k$  principal components. We can expect  $k$  to be small when the variables are strongly correlated. When the variables are mutually uncorrelated, principal components analysis is not useful.

Here is a summary of the properties of principal components.

**Lemma 12.9.1.** *Let  $Z_s = \mathbf{u}_s^T (\mathbf{Y} - \boldsymbol{\mu})$ ,  $s = 1, 2, \dots$  Then:*

- (a)  $E[Z_s] = 0$ .
- (b)  $\text{Var}[Z_s] = \lambda_s$ .
- (c)  $\text{Cov}[Z_s, Z_{s'}] = 0$  if  $s \neq s'$ .
- (d)  $\text{Corr}[Y_l, Z_s] = \sqrt{\lambda_s} u_{ls} / \sqrt{\sigma_{ll}}$ , where  $Y_l$  is the  $l$ th component of  $\mathbf{Y}$ ,  $u_{ls}$  is the  $l$ th element of the vector  $\mathbf{u}_s$  and  $\sigma_{ll}$  is the  $l$ th diagonal element of  $\Sigma$ .

### 12.9.1 Regression Interpretation of Principal Components

Let  $Z_1 = \mathbf{e}_1^T(\mathbf{Y} - \mu), \dots, Z_k = \mathbf{e}_k^T(\mathbf{Y} - \mu)$  be  $k$  linear functions of the random vector  $\mathbf{Y}$  such that the random variables  $Z_1, \dots, Z_k$  are mutually uncorrelated. Now consider the problem of predicting  $Y_j$ , the  $j$ th component of  $\mathbf{Y}$ , from  $Z_1, \dots, Z_k$  using a linear regression and let  $\tau_j^2$  be the prediction error. Then  $\sum_{j=1}^p \tau_j^2$  is the total prediction error for predicting  $Y_1, \dots, Y_p$  (each separately) using  $Z_1, \dots, Z_k$ . If  $\sum_{j=1}^p \tau_j^2$  is quite small, then we may conclude that the information in the vector  $\mathbf{Y}$  can be well summarized by  $Z_1, \dots, Z_k$ . It turns out that  $\sum_{j=1}^p \tau_j^2$  is minimized when  $Z_1, \dots, Z_k$  are the first  $k$  principal components of  $\mathbf{Y}$ .

Let us now justify this regression interpretation. Note that  $Y_j - \mu_j, Z_1, \dots, Z_k$  have zero means, and hence the intercept term for the regression of  $Y_j - \mu_j$  on  $Z_1, \dots, Z_k$  is zero. Clearly,

$$\tau_j^2 = \min_{\beta_{j1}, \dots, \beta_{jk}} E \left\{ Y_j - \mu_j - \sum_{s=1}^k \beta_{js} Z_s \right\}^2.$$

Since  $Z_1, \dots, Z_k$  are assumed to be mutually uncorrelated, the solution to this minimization problem is given by  $\beta_{js}^* = \text{Cov}[Y_j, Z_s]/\text{Var}[Z_s]$ ,  $1 \leq s \leq k$ , and consequently

$$\tau_j^2 = \text{Var}[Y_j] - \sum_{s=1}^k \text{Cov}[Y_j, Z_s]^2 / \text{Var}[Z_s].$$

Let  $\boldsymbol{\sigma}_1, \dots, \boldsymbol{\sigma}_p$  be the column vectors of the covariance matrix  $\Sigma$ , that is,  $\Sigma = [\boldsymbol{\sigma}_1, \dots, \boldsymbol{\sigma}_p]$ . Then  $\text{Cov}[Y_j, Z_s] = \mathbf{e}_s^T \boldsymbol{\sigma}_j$ ,  $\text{Var}[Z_s] = \mathbf{e}_s^T \Sigma \mathbf{e}_s$ , and we have

$$\tau_j^2 = \text{Var}[Y_j] - \sum_{s=1}^k \frac{(\mathbf{e}_s^T \boldsymbol{\sigma}_j)^2}{\mathbf{e}_s^T \Sigma \mathbf{e}_s}.$$

Hence

$$\begin{aligned} \sum_{j=1}^p \tau_j^2 &= \sum_{j=1}^p \text{Var}[Y_j] - \sum_{j=1}^p \sum_{s=1}^k \frac{(\mathbf{e}_s^T \boldsymbol{\sigma}_j)^2}{\mathbf{e}_s^T \Sigma \mathbf{e}_s} \\ &= \sum_{j=1}^p \text{Var}[Y_j] - \sum_{s=1}^p \sum_{j=1}^k \frac{(\mathbf{e}_s^T \boldsymbol{\sigma}_j)^2}{\mathbf{e}_s^T \Sigma \mathbf{e}_s}. \end{aligned}$$

Since  $\sum_{j=1}^p \text{Var}[Y_j] = \text{trace}(\boldsymbol{\Sigma})$  and

$$\sum_{j=1}^p (\mathbf{e}_s^T \boldsymbol{\sigma}_j)^2 = \mathbf{e}_s^T \sum_{j=1}^p \boldsymbol{\sigma}_j \boldsymbol{\sigma}_j^T \mathbf{e}_s = \mathbf{e}_s^T \boldsymbol{\Sigma}^2 \mathbf{e}_s,$$

we have

$$\sum_{j=1}^p \tau_j^2 = \text{trace}(\boldsymbol{\Sigma}) - \sum_{s=1}^k \frac{\mathbf{e}_s^T \boldsymbol{\Sigma}^2 \mathbf{e}_s}{\mathbf{e}_s^T \boldsymbol{\Sigma} \mathbf{e}_s}.$$

Minimizing  $\sum_{j=1}^p \tau_j^2$  with respect to  $\mathbf{e}_1, \dots, \mathbf{e}_k$  subject to the constraints  $\mathbf{e}_s^T \boldsymbol{\Sigma} \mathbf{e}_{s'} = \text{Cov}[Z_s, Z_{s'}] = 0$ ,  $1 \leq s \neq s' \leq k$ , is equivalent to maximizing

$$\sum_{s=1}^k \frac{\mathbf{e}_s^T \boldsymbol{\Sigma}^2 \mathbf{e}_s}{\mathbf{e}_s^T \boldsymbol{\Sigma} \mathbf{e}_s} \quad (2)$$

with respect to  $\mathbf{e}_1, \dots, \mathbf{e}_k$  subject to the constraints  $\mathbf{e}_s^T \boldsymbol{\Sigma} \mathbf{e}_{s'} = 0$ ,  $1 \leq s \neq s' \leq k$ . This is a generalized eigenvalue problem described in Section B.4. If  $\lambda_1 \geq \lambda_2 \geq \dots$  are the generalized eigenvalues of  $\mathbf{A} = \boldsymbol{\Sigma}^2$  with respect to  $\mathbf{B} = \boldsymbol{\Sigma}$ , the maximum value of Eq. (2) is given by  $\lambda_1 + \dots + \lambda_k$ . Since  $\lambda_1, \lambda_2, \dots$  are the eigenvalues of  $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2} = \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}^2 \boldsymbol{\Sigma}^{-1/2} = \boldsymbol{\Sigma}$  with the corresponding orthonormal eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots$ , the maximum of Eq. (2) occurs at  $\mathbf{e}_s = \mathbf{u}_s$ ,  $s = 1, \dots, k$ .

This proves that among all the  $k$  mutually uncorrelated linear functions of  $\mathbf{Y}$ , the principal components  $Z_1 = \mathbf{u}_1^T (\mathbf{Y} - \boldsymbol{\mu}), \dots, Z_k = \mathbf{u}_k^T (\mathbf{Y} - \boldsymbol{\mu})$  are the best linear predictors of  $Y_1, \dots, Y_p$ .

### 12.9.2 Estimation of Principal Components

Till now we have discussed the concept of principal components for a random vector from a theoretical standpoint. We now take up the issue of estimating them from the data. Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . Note that we have made no assumption on the distribution of  $\mathbf{Y}_i$  (such as normality).

Even though the covariance matrix of the standardized variables is the correlation matrix, for notational simplicity, we use  $\boldsymbol{\Sigma}$  to denote both the covariance and correlation matrices. Similarly we use  $\mathbf{S}$  to denote both the sample covariance and correlation matrices. Note that if we want to find the principal components of  $\mathbf{Y} - \boldsymbol{\mu}$ , then we deal with the covariance matrix of  $\mathbf{Y}$ . Whereas if we want to find the principal components of  $[\text{diag}(\boldsymbol{\Sigma})]^{-1/2} (\mathbf{Y} - \boldsymbol{\mu})$  (ie, the standardized variables), we deal with the correlation matrix.

As before, let  $\lambda_1 \geq \lambda_2 \geq \dots$  be the eigenvalues of  $\boldsymbol{\Sigma}$  with  $\mathbf{u}_1, \mathbf{u}_2, \dots$  as the corresponding orthonormal eigenvectors. Similarly, let  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  be the eigenvalues of  $\mathbf{S}$  with the corresponding orthonormal eigenvectors  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$ . It is easy to guess that  $\hat{\lambda}_j$  estimates  $\lambda_j$

and  $\hat{\mathbf{u}}_j$  estimates  $\pm \mathbf{u}_j$ ,  $j = 1, \dots, p$ . The reason for “ $\pm$ ” is that if  $\mathbf{u}_j$  is an eigenvector of  $\Sigma$  with the eigenvalue  $\lambda_j$ , then  $-\mathbf{u}_j$  is also an eigenvector of  $\Sigma$  with the same eigenvalue  $\lambda_j$ . Estimated principal components are  $\hat{Z}_s = \hat{\mathbf{u}}_s^T (\mathbf{Y} - \bar{\mathbf{Y}})$ ,  $s = 1, \dots, p$ . Recall that the population (or theoretical) principal components have zero means and are uncorrelated. The estimated (sample) principal components have similar properties. Define the scores of the  $s$ th principal component to be  $\hat{Z}_{is} = \hat{\mathbf{u}}_s^T (\mathbf{Y}_i - \bar{\mathbf{Y}})$ ,  $i = 1, \dots, n$ . A result analogous to Lemma 12.9.1 holds for the sample principal components when  $E[\cdot]$ ,  $\text{Var}[\cdot]$ , etc., are replaced by the sample mean, sample variance, etc., and this is left as an exercise. In particular,

$$\begin{aligned}\sum_{i=1}^n \hat{Z}_{is}/n &= 0, \quad \sum_{i=1}^n \hat{Z}_{is}^2/(n-1) = \hat{\lambda}_s, \text{ and} \\ \sum_{i=1}^n \hat{Z}_{is} \hat{Z}_{is'}/(n-1) &= 0, \quad s \neq s'.\end{aligned}$$

### Scree Plot

The ratio  $\hat{\lambda}_k / \sum \hat{\lambda}_j = \hat{\lambda}_k / \text{trace}(\mathbf{S})$  estimates the proportion of total variability of  $\mathbf{Y}$  explained by the  $k$ th principal component. For this reason, it is useful to plot  $\hat{\lambda}_k / \text{trace}(\mathbf{S})$  against  $k$ . This plot graphically displays how the eigenvalues decay. It is also useful to plot the cumulative ratio  $\hat{\pi}_k = (\hat{\lambda}_1 + \dots + \hat{\lambda}_k) / (\hat{\lambda}_1 + \dots + \hat{\lambda}_p)$  against  $k$ . This gives us an idea of how much dimensionality reduction is possible since  $\hat{\pi}_k$  is an estimate of  $\pi_k = (\lambda_1 + \dots + \lambda_k) / (\lambda_1 + \dots + \lambda_p)$ , the proportion of the total variability of  $\mathbf{Y}$  explained by the first  $k$  principal components.

### 12.9.3 Asymptotic Results in Principal Components Analysis

Asymptotic distributions of sample eigenvalues and sample eigenvectors are somewhat complicated. It is important to point out that the asymptotic distributions depend on the distribution of the multivariate population from which the observations are taken. This is quite unlike the limit theorems for the estimated mean vectors and associated statistics. Here we only deal with the case when the eigenvalues of the population covariance matrix  $\Sigma$  are distinct. When some of the eigenvalues have multiplicities larger than 1, the asymptotic distributions of the sample eigenvalues and sample eigenvectors are not even normal (they are mixtures of normals). We do not discuss such cases here.

As before, let the spectral decompositions of  $\Sigma$  and sample covariance matrix  $\mathbf{S}$  be

$$\Sigma = \sum_{j=1}^p \lambda_j \mathbf{u}_j \mathbf{u}_j^T \text{ and } \mathbf{S} = \sum_{j=1}^p \hat{\lambda}_j \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T,$$

where  $\lambda_1 \geq \lambda_2 \geq \dots$  and  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$ . We assume all the eigenvalues are distinct (ie,  $\lambda_1 > \lambda_2 > \dots > \lambda_p$ ). For continuous distributions, it can be shown that the sample eigenvalues  $\hat{\lambda}_1, \dots, \hat{\lambda}_p$  are distinct with probability 1.

The  $j$ th principal component of  $\mathbf{Y}$  is  $Z_j = \mathbf{u}_j^T(\mathbf{Y} - \boldsymbol{\mu})$ . We know that  $E[Z_j] = 0$  and  $\text{Var}[Z_j] = \lambda_j$ . We also know that  $Z_1, \dots, Z_p$  are mutually uncorrelated. If the population of  $\mathbf{Y}$  is normal then  $Z_1, \dots, Z_p$  are independent and  $Z_j \sim N(0, \lambda_j)$ ,  $j = 1, \dots, p$ .

We will write down the asymptotic distributions of the sample eigenvalues and sample eigenvectors. These results can be proved using the perturbation theory of matrices which is outside the scope of this book.

**Theorem 12.9.1.** *Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , and assume that the eigenvalues of  $\boldsymbol{\Sigma}$  are distinct. All the results below are true assuming that  $n \rightarrow \infty$ .*

- (a)  $\sqrt{n}(\hat{\lambda} - \lambda) \xrightarrow{\mathcal{D}} N_p(\mathbf{0}, \mathbf{W})$ , where element  $(j, k)$  of  $\mathbf{W}$  is given by  $\text{Cov}[Z_j^2, Z_k^2]$ , where  $\lambda$  and  $\hat{\lambda}$  are  $p$ -dim vectors of the eigenvalues  $\lambda_1, \dots, \lambda_p$  and their estimates.
- (b) If the population is normal (ie,  $\{\mathbf{Y}_i\}$  are iid  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ), then the matrix  $\mathbf{W}$  in part (a) is a diagonal matrix whose  $j$ th diagonal element is given by  $2\lambda_j^2$ .
- (c) For any  $1 \leq i \leq k$ ,  $\sqrt{n}(\hat{\mathbf{u}}_i \pm \mathbf{u}_i) \xrightarrow{\mathcal{D}} N_p(\mathbf{0}, \mathbf{R}_i)$ , where

$$\mathbf{R}_i = \sum_{j \neq i} \sum_{k \neq i} \frac{\text{Cov}[Z_j Z_i, Z_k Z_i]}{(\lambda_j - \lambda_i)(\lambda_k - \lambda_i)} \mathbf{u}_j \mathbf{u}_k^T.$$

- (d) If the population is normal, then the matrix  $\mathbf{R}_i$  in part (c) has the simplified form

$$\mathbf{R}_i = \lambda_i \sum_{j \neq i} \frac{\lambda_j}{(\lambda_j - \lambda_i)^2} \mathbf{u}_j \mathbf{u}_j^T.$$

The results given above allow us to construct confidence intervals for the eigenvalues.

*Remark 12.9.2.* Even though  $\{Z_j\}$  are mutually uncorrelated, the same is not necessarily true for  $\{Z_j^2\}$ . If the population is normal,  $\{Z_j^2\}$  are mutually independent, and hence mutually uncorrelated. Thus, under normality,  $\text{Cov}[Z_j^2, Z_k^2] = 0$ ,  $j \neq k$ , and  $\text{Cov}[Z_j^2, Z_k^2] = \text{Var}[Z_j^2] = 2\lambda_j^2$  when  $j = k$ . Therefore, part (b) of the last result follows from part (a). When  $j \neq i$  and  $k \neq i$ ,

$$\text{Cov}[Z_j Z_i, Z_k Z_i] = E[Z_j Z_k Z_i^2] - E[Z_j Z_i]E[Z_k Z_i] = E[Z_j Z_k Z_i^2].$$

Under normality, when  $j$ ,  $k$ , and  $i$  are all distinct,  $E[Z_j Z_k Z_i^2] = 0$ . When  $j = k \neq i$ , then  $E[Z_j Z_k Z_i^2] = E[Z_j^2]E[Z_i^2] = \lambda_j \lambda_i$ . This shows that part (d) of the last theorem follows from part (c).

When the population is normal, estimation of  $\mathbf{W}$  or  $\mathbf{R}_i$  is rather easy since  $\hat{\mathbf{W}}$  is a diagonal matrix with diagonal elements  $2\hat{\lambda}_1^2, \dots, 2\hat{\lambda}_p^2$ , and

$$\hat{\mathbf{R}}_i = \hat{\lambda}_i \sum_{j \neq i} \frac{\hat{\lambda}_j}{(\hat{\lambda}_j - \hat{\lambda}_i)^2} \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T.$$

However, if the population is not normal, then we need to estimate these matrices ( $\mathbf{W}$  and  $\mathbf{R}_i$ ) using the principal component scores. For the  $i$ th principal component, the scores are  $\{\hat{Z}_{ti} = \mathbf{u}_i^T(\mathbf{Y}_t - \bar{\mathbf{Y}}) : t = 1, \dots, n\}$ . Estimate of element  $(j, k)$  of  $\mathbf{W}$  is

$$\begin{aligned}\hat{W}_{jk} &= \text{sample covariance of } \{(\hat{Z}_{tj}^2, \hat{Z}_{tk}^2): t = 1, \dots, n\} \\ &= (n-1)^{-1} \left\{ \sum_{t=1}^n \hat{Z}_{tj}^2 \hat{Z}_{tk}^2 - n \hat{\lambda}_j \hat{\lambda}_k \right\}.\end{aligned}$$

Similarly

$$\hat{\mathbf{R}}_i = \sum_{j \neq i} \sum_{k \neq i} \frac{\nu(j, k, i)}{(\hat{\lambda}_j - \hat{\lambda}_i)(\hat{\lambda}_k - \hat{\lambda}_i)} \hat{\mathbf{u}}_j \hat{\mathbf{u}}_k^T,$$

where  $\nu(j, k, i)$  is the sample estimate of  $\text{Cov}[Z_j Z_i, Z_k Z_i]$  and is given by

$$\nu(j, k, i) = (n-1)^{-1} \sum_{t=1}^n \hat{Z}_{tj} \hat{Z}_{tk} \hat{Z}_{ti}^2.$$

### *Confidence Interval for $\lambda_j$*

An approximate confidence interval for  $\lambda_j$  with a confidence coefficient  $1 - \alpha$  is given by  $\hat{\lambda}_j \pm z_{\alpha/2} \sqrt{\hat{W}_{jj}/n}$ , where  $z_{\alpha/2}$  is the  $(1 - \alpha/2)$ -quantile of the standard normal distribution. Note that  $\hat{W}_{jj} = 2\hat{\lambda}_j^2$  for the normal case, whereas  $\hat{W}_{jj}$  is the sample variance of  $\{\hat{Z}_{tj}^2: t = 1, \dots, n\}$  in the general case. Sometimes it may be preferable to construct a confidence interval for  $\log \lambda_j$ . By the delta method,

$$\sqrt{n}(\log \hat{\lambda}_j - \log \lambda_j) \xrightarrow{\mathcal{D}} N(0, W_{jj}/\lambda_j^2).$$

If the population is normal, then the natural logarithm is a variance stabilizing transformation since  $W_{jj}/\lambda_j^2 = 2$ , and in such a case, an approximate confidence interval for  $\log \lambda_j$  with confidence coefficient  $1 - \alpha$  is given by  $\log \hat{\lambda}_j \pm z_{\alpha/2} \sqrt{2/n}$ . In the general case, an approximate confidence interval for  $\log \lambda_j$  is  $\log \hat{\lambda}_j \pm z_{\alpha/2} \sqrt{\hat{W}_{jj}/(n\hat{\lambda}_j^2)}$ .

### *Estimation of $\pi_k = (\lambda_1 + \dots + \lambda_k)/(\lambda_1 + \dots + \lambda_p)$*

Recall that  $\pi_k$  is the proportion of variability of  $\mathbf{Y}$  explained by the first  $k$  principal components. An estimate of  $\pi_k$  is

$$\begin{aligned}\hat{\pi}_k &= (\hat{\lambda}_1 + \dots + \hat{\lambda}_k)/(\hat{\lambda}_1 + \dots + \hat{\lambda}_p) \\ &= (\hat{\lambda}_1 + \dots + \hat{\lambda}_k)/\text{trace}(\mathbf{S}).\end{aligned}$$

In order to obtain a confidence interval for  $\pi_k$ , we need to find the asymptotic distribution of  $\hat{\pi}_k$ . Let

$$g(\lambda) = (\lambda_1 + \dots + \lambda_k)/(\lambda_1 + \dots + \lambda_p)$$

and let  $\mathbf{g}_1(\lambda)$  be the vector of partial derivatives of  $g$ . By the delta method,

$$\sqrt{n}(\hat{\pi}_k - \pi_k) \xrightarrow{\mathcal{D}} N(0, \mathbf{g}_1(\lambda)^T \mathbf{W} \mathbf{g}_1(\lambda)) \text{ as } n \rightarrow \infty.$$

## 12.10 Factor Analysis

Let  $\mathbf{Y}$  be a  $p$ -dim vector with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . A factor model with  $k$  common factors ( $k \leq p$ ) is

$$\mathbf{Y} - \boldsymbol{\mu} = \mathbf{L}\mathbf{f} + \boldsymbol{\epsilon} \text{ or } \mathbf{Y} = \boldsymbol{\mu} + \mathbf{L}\mathbf{f} + \boldsymbol{\epsilon}, \quad (3)$$

where  $\mathbf{L}$  is a  $p \times k$  matrix of factor loadings (nonrandom),  $\mathbf{f}$  is a  $k \times 1$  vector of common factors (random), and  $\boldsymbol{\epsilon}$  is a  $p \times 1$  vector of specific factors (random). We assume that

$$\mathbb{E}[\mathbf{f}] = \mathbf{0}, \text{ Cov}[\mathbf{f}] = \mathbf{I}, \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \text{ Cov}[\boldsymbol{\epsilon}] = \boldsymbol{\Psi}, \text{ and } \text{Cov}[\mathbf{f}, \boldsymbol{\epsilon}] = \mathbf{0}.$$

Here  $\boldsymbol{\Psi}$  is assumed to be a diagonal matrix with positive diagonal elements  $\psi_1, \dots, \psi_p$ .

Note that if the factor model is correct, then  $\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}$  and  $\text{Cov}[\mathbf{Y}] = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$ . The goal of factor analysis is to approximate the covariance matrix  $\boldsymbol{\Sigma}$  by a matrix of the form  $\mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$ , where  $\mathbf{L}$  is  $p \times k$  and  $\boldsymbol{\Psi}$  is a diagonal matrix with positive diagonal elements. In other words, factor analysis can be viewed as a dimensionality reduction method, and, in order for this method to be useful,  $k$  should be as small as possible (at least in comparison to  $p$ ).

Let  $Y_i$  be the  $i$ th component of  $\mathbf{Y}$ . If the factor model is correct, then  $\mathbb{E}[Y_i] = \mu_i$  and

$$\text{Var}[Y_i] = l_{i1}^2 + \dots + l_{ik}^2 + \psi_i = h_i^2 + \psi_i, \quad i = 1, \dots, p,$$

where  $h_i^2 = l_{i1}^2 + \dots + l_{ik}^2$  is called the “communality” and  $\psi_i$  is called the specific variance.

It is important to note that the factor model in Eq. (3) is not identifiable. If  $\mathbf{G}$  is a  $k \times k$  orthogonal matrix and  $\tilde{\mathbf{L}} = \mathbf{L}\mathbf{G}$ , then we can write

$$\mathbf{Y} - \boldsymbol{\mu} = \mathbf{L}\mathbf{f} + \boldsymbol{\epsilon} = \tilde{\mathbf{L}}\tilde{\mathbf{f}} + \boldsymbol{\epsilon},$$

where  $\tilde{\mathbf{f}} = \mathbf{G}^T \mathbf{f}$ . Note that  $\mathbb{E}[\tilde{\mathbf{f}}] = \mathbf{0}$  and  $\text{Cov}[\tilde{\mathbf{f}}] = \mathbf{I}$ , and

$$\text{Cov}[\mathbf{Y}] = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi} = \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T + \boldsymbol{\Psi}.$$

### 12.10.1 Estimation of $L$ and $\boldsymbol{\Psi}$

Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be a random sample from a population with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . As usual the mean vector is estimated by the sample mean. The goal is to estimate  $\mathbf{L}$  and  $\boldsymbol{\Psi}$  from the data assuming that the factor model given in Eq. (3) is appropriate (ie,  $\text{Cov}[\mathbf{Y}] = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$ ). The following two methods are widely used for estimating  $\mathbf{L}$  and  $\boldsymbol{\Psi}$ :

- (a) principal components,
- (b) maximum likelihood.

### Principal Factor Analysis

Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$  be the eigenvalues of the sample covariance matrix  $\mathbf{S}$  with the corresponding orthonormal eigenvectors  $\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_p$ . Then the  $p \times k$  matrix

$\hat{\mathbf{L}} = [\hat{\lambda}_1^{1/2} \hat{\mathbf{u}}_1, \dots, \hat{\lambda}_k^{1/2} \hat{\mathbf{u}}_k]$  is taken to be an estimate of  $\mathbf{L}$ , that is, the columns of  $\hat{\mathbf{L}}$  are  $\hat{\mathbf{l}}_j = \hat{\lambda}_j^{1/2} \hat{\mathbf{u}}_j, j = 1, \dots, k$ , and hence  $\hat{\mathbf{L}}\hat{\mathbf{L}}^T = \sum_{j=1}^k \hat{\lambda}_j \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T$ . Here

$$\hat{h}_i^2 = \hat{l}_{i1}^2 + \dots + \hat{l}_{ik}^2, \quad \hat{\psi}_i = s_{ii} - \hat{h}_i^2,$$

where  $s_{ii}$  is the  $i$ th diagonal element of the sample covariance matrix  $\mathbf{S}$ . An estimate of the proportion of total variability of  $\mathbf{Y}$  explained by the first factor is  $\|\hat{\mathbf{l}}_1\|^2/\text{trace}(\mathbf{S}) = \hat{\lambda}_1/\text{trace}(\mathbf{S})$ . In general, an estimate of the proportion of total variability of  $\mathbf{Y}$  explained by the  $j$ th factor is  $\hat{\lambda}_j/\text{trace}(\mathbf{S}), j = 1, \dots, k$ .

If the sample correlation matrix is used in the analysis instead of the sample covariance matrix, then  $\hat{\mathbf{L}} = [\hat{\lambda}_1^{1/2} \hat{\mathbf{u}}_1, \dots, \hat{\lambda}_k^{1/2} \hat{\mathbf{u}}_k]$ , where  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  are the eigenvalues of the sample correlation matrix with the corresponding normalized eigenvectors  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$ . In this case,  $\hat{h}_i^2 = \hat{l}_{i1}^2 + \dots + \hat{l}_{ik}^2$ ,  $\hat{\psi}_i = 1 - \hat{h}_i^2$ , and the proportion of the total variability explained by the  $j$ th factor is  $\hat{\lambda}_j/p, j = 1, 2, \dots$

## Maximum Likelihood

This method typically assumes that  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  are iid  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The maximum likelihood estimates of  $\mathbf{L}$  and  $\boldsymbol{\Psi}$  are obtained by maximizing the likelihood under the constraint that  $\mathbf{L}\boldsymbol{\Psi}^{-1}\mathbf{L}^T$  is a diagonal matrix. This constraint eliminates the problem of nonuniqueness of the factor loading matrix  $\mathbf{L}$ . Explicit expression for the estimates are not available. However, computer packages such R or MATLAB can be used to obtain the estimates of  $\mathbf{L}$  and  $\boldsymbol{\Psi}$ . Estimate of the proportion of the total variability explained by the  $j$ th factor is  $\|\hat{\mathbf{l}}_j\|^2/\text{trace}(\mathbf{S}), j = 1, \dots, k$ , where  $\hat{\mathbf{l}}_1, \dots, \hat{\mathbf{l}}_k$  are the columns of  $\hat{\mathbf{L}}$ .

If the correlation matrix is used in the analysis, the estimate of the proportion of the total variability explained by the  $j$ th factor is  $\|\hat{\mathbf{l}}_j\|^2/p, j = 1, \dots, k$ .

## How Many Factors

Since we are trying to estimate  $\boldsymbol{\Sigma}$  by a matrix of the form  $\mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$ , a reasonable way to determine  $k$  is to examine the residual matrix  $\mathbf{S} - (\hat{\mathbf{L}}\hat{\mathbf{L}}^T + \hat{\boldsymbol{\Psi}})$ . If all the elements of the residual matrix are small, then we may assume that  $\mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$  is a good approximation to  $\boldsymbol{\Sigma}$ . A more formal way to decide the number of factors is to minimize an AIC-type criterion over  $k$ . However, we do not address that issue here.

## Factor Rotation

If  $\hat{\mathbf{L}}$  is an estimate of the factor loading matrix  $\mathbf{L}$  and  $\hat{\mathbf{L}}_* = \hat{\mathbf{L}}\mathbf{G}$ , where  $\mathbf{G}$  is a  $k \times k$  orthogonal matrix, then  $\hat{\mathbf{L}}\hat{\mathbf{L}}^T + \hat{\boldsymbol{\Psi}} = \hat{\mathbf{L}}_*\hat{\mathbf{L}}_*^T + \hat{\boldsymbol{\Psi}}$ . So we may take  $\hat{\mathbf{L}}_*$  to be also a valid estimate of  $\mathbf{L}$ . It has been suggested that one should choose  $\hat{\mathbf{L}}_*$  so that for each column of  $\hat{\mathbf{L}}_*$ , some of the elements are relatively large and the rest are small. Kaizer's Varimax Rotation tries to achieve this. The procedure is described below.

Let  $\tilde{l}_{*ij} = \hat{l}_{*ij}/\hat{h}_i$ ,  $j = 1, \dots, k$ ,  $i = 1, \dots, p$ , where  $\hat{l}_{*ij}$  is the  $j^{th}$  element in the  $i^{th}$  row of  $\hat{\mathbf{L}}_*$ . Now consider the  $p \times k$  matrix  $\tilde{\mathbf{L}}_*$  whose elements are given by  $\tilde{l}_{*ij}$ 's. Note that each row of  $\tilde{\mathbf{L}}_*$  has length 1 since for any  $1 \leq i \leq p$ ,  $\sum_{j=1}^k \tilde{l}_{*ij}^2 = 1$ . We consider a criterion which maximizes the variability of the squares of the elements in each column of  $\tilde{\mathbf{L}}_*$ . For  $1 \leq j \leq k$ , consider a measure of the variability of  $\{\tilde{l}_{*1j}^2, \dots, \tilde{l}_{*pj}^2\}$

$$(1/p) \sum_{i=1}^p \left\{ \tilde{l}_{*ij}^4 - \left( \sum_i \tilde{l}_{*ij}^2 / p \right)^2 \right\}.$$

Now add these measures over  $j = 1, \dots, k$ , to obtain the following criterion

$$\begin{aligned} Q &= \sum_{j=1}^k (1/p) \sum_{i=1}^p \left\{ \tilde{l}_{*ij}^4 - \left( \sum_i \tilde{l}_{*ij}^2 / p \right)^2 \right\} \\ &= \sum_{j=1}^k [\text{variance of } \{\tilde{l}_{*ij}^2 : i = 1, \dots, p\}]. \end{aligned}$$

We can maximize  $Q$  in order to obtain  $\hat{\mathbf{L}}_*$  which is a rotated version of  $\hat{\mathbf{L}}$ . There is no explicit expression for the rotated matrix, but one can obtain this estimate by using a computer package.

### 12.10.2 Prediction of Common Factors

Prediction of common factors in a factor analysis setting is similar to prediction in a random- or mixed-effect model discussed in Section 11.10 of Chapter 11. Assume that the factor model in Eq. (3) holds,  $\mu$ ,  $\mathbf{L}$ , and  $\Psi$  are known (or estimated using past observations), then the goal is to predict  $\mathbf{f}$  when  $\mathbf{Y}$  is observed. We restrict ourselves to predictors which are linear functions of  $\mathbf{Y}$ . A predictor  $\hat{\mathbf{f}}$  of  $\mathbf{f}$  is called unbiased if  $E[\hat{\mathbf{f}} - \mathbf{f}] = \mathbf{0}$ . In the following discussion we assume that  $\mathbf{L}$  has rank  $k$ , and the diagonal elements of  $\Psi$  are positive. Since we can center  $\mathbf{Y}$  by subtracting the mean  $\mu$ , we assume that  $\mu = \mathbf{0}$  in the subsequent discussion.

If we ignore that  $\mathbf{f}$  is random and minimize the ordinary least squares criterion  $\|\mathbf{Y} - \mathbf{Lf}\|$  with respect to  $\mathbf{f}$ , we get a predictor of the form  $\hat{\mathbf{f}}^{(1)} = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{Y}$ . Similarly, we may try to obtain a weighted least squares predictor of  $\mathbf{f}$  as follows. If we premultiply both sides of Eq. (3) by  $\Psi^{-1/2}$ , we get model of the form  $\tilde{\mathbf{Y}} = \mathbf{Lf} + \tilde{\boldsymbol{\epsilon}}$ , where  $\tilde{\mathbf{Y}} = \Psi^{-1/2} \mathbf{Y}$ ,  $\tilde{\mathbf{L}} = \Psi^{-1/2} \mathbf{L}$ , and  $\tilde{\boldsymbol{\epsilon}} = \Psi^{-1/2} \boldsymbol{\epsilon}$ . Noting that  $E[\tilde{\boldsymbol{\epsilon}}] = \mathbf{0}$  and  $\text{Cov}[\tilde{\boldsymbol{\epsilon}}] = \mathbf{I}$ , we may minimize  $\|\tilde{\mathbf{Y}} - \tilde{\mathbf{L}}\mathbf{f}\|^2 = (\mathbf{Y} - \mathbf{Lf})^T \Psi^{-1} (\mathbf{Y} - \mathbf{Lf})$  with respect to  $\mathbf{f}$ , and this leads to another linear predictor  $\hat{\mathbf{f}}^{(2)} = (\mathbf{L}^T \Psi^{-1} \mathbf{L})^{-1} \mathbf{L}^T \Psi^{-1} \mathbf{Y}$ . Both  $\hat{\mathbf{f}}^{(1)}$  and  $\hat{\mathbf{f}}^{(2)}$  are unbiased linear predictors of  $\mathbf{f}$ , but they are not the best in terms of prediction error.

As in [Section 11.10](#) of [Chapter 11](#), we may define the concept of best linear unbiased predictor as follows.

**Definition 12.10.1.** A linear function  $\hat{\mathbf{f}}$  of  $\mathbf{Y}$  is called a best linear unbiased predictor of  $\mathbf{f}$  if

- (i)  $\hat{\mathbf{f}}$  is an unbiased predictor of  $\mathbf{f}$ , that is,  $E[\hat{\mathbf{f}} - \mathbf{f}] = \mathbf{0}$ .
- (ii) For any  $\mathbf{a} \in \mathbb{R}^k$ ,  $E[\mathbf{a}^T \hat{\mathbf{f}} - \mathbf{a}^T \mathbf{f}]^2 \leq E[\mathbf{l}^T \mathbf{Y} - \mathbf{a}^T \mathbf{f}]^2$  for all linear unbiased predictors  $\mathbf{l}^T \mathbf{Y}$  of  $\mathbf{a}^T \mathbf{f}$ ,  $\mathbf{l} \in \mathbb{R}^p$ .

Arguments used in [Sections 11.10.1](#) and [11.10.2](#) lead to the following best linear predictor

$$\hat{\mathbf{f}}^{(3)} = (\mathbf{L}^T \boldsymbol{\Psi}^{-1} \mathbf{L} + \mathbf{I})^{-1} \mathbf{L}^T \boldsymbol{\Psi}^{-1} \mathbf{Y}.$$

We leave it to the reader to prove that the best linear unbiased predictor of  $\mathbf{f}$  is unique and is equal to  $\hat{\mathbf{f}}^{(3)}$  as given above.

## 12.11 Classification and Discrimination

Suppose that we have  $k$  populations (each  $p$ -dim) with means  $\mu_1, \dots, \mu_k$  and the same covariance matrix  $\Sigma$ . If we have an observation vector  $\mathbf{y}$  from one of these populations, then the goal of the classification problem is to guess which population  $\mathbf{y}$  comes from. If  $\mathbf{y}$  is closer to  $\mu_i$  than all other means, then a reasonable guess is that  $\mathbf{y}$  comes from population  $i$ . It turns out that this intuition is also mathematically valid. Recall that Mahalanobis distance between  $\mathbf{y}$  and  $\mu_i$  is  $\Delta^2(\mathbf{y}, \mu_i) = (\mathbf{y} - \mu_i)^T \Sigma^{-1} (\mathbf{y} - \mu_i)$ . So a reasonable rule is: allocate  $\mathbf{y}$  to population  $i$  if

$$\begin{aligned} \Delta^2(\mathbf{y}, \mu_i) &< \Delta^2(\mathbf{y}, \mu_j), \quad \text{for all } j \neq i, \text{ ie,} \\ -2\mu_i^T \Sigma^{-1} \mathbf{y} + \mu_i^T \Sigma^{-1} \mu_i &< -2\mu_j^T \Sigma^{-1} \mathbf{y} + \mu_j^T \Sigma^{-1} \mu_j, \quad \text{for all } j \neq i. \end{aligned} \tag{4}$$

This is called a linear discriminant rule since the criterion for discrimination between the populations depends linearly on  $\mathbf{y}$ .

In some cases, the prior probabilities  $\{\pi_1, \dots, \pi_k\}$  of the populations need to be taken into account. For such a case, the rule given above is modified: now we allocate  $\mathbf{y}$  to population  $i$  if

$$\Delta^2(\mathbf{y}, \mu_i) - 2 \log \pi_i < \Delta^2(\mathbf{y}, \mu_j) - 2 \log \pi_j, \quad \text{for all } j \neq i. \tag{5}$$

Note that the rule defined by the inequalities given in (5) is also a linear discriminant rule, and the rule in (4) is a special case of the rule in (5) when the prior is noninformative, that is, when  $\pi_1 = \pi_2 = \dots = \pi_k = 1/k$ .

**Note.** If the  $k$  populations have different covariance matrices, then the rules defined by the inequalities in (4) and (5) need to be modified, and, in such a case, we are led to what is known as a quadratic discriminant rule. We discuss this issue below.

### 12.11.1 Bayes' Rule for Classification

Let us assume that the chance that the random vector  $\mathbf{Y}$  is from population  $i$  (with pdf  $f_i$ ) is  $\pi_i$ ,  $i = 1, \dots, k$ , where  $\pi_1 + \dots + \pi_k = 1$ . Let  $J$  be a discrete variable taking values  $1, \dots, k$ , with  $P[J = i] = \pi_i$ . In this framework, the conditional pdf of  $\mathbf{Y}$  given  $J = i$  is  $f_i$ . The marginal pdf of  $\mathbf{Y}$  is  $f(\mathbf{y}) = \pi_1 f_1(\mathbf{y}) + \dots + \pi_k f_k(\mathbf{y})$ , and the conditional probability of  $J = i$  given  $\mathbf{Y} = \mathbf{y}$  is  $P[J = i | \mathbf{Y} = \mathbf{y}] = \pi_i f_i(\mathbf{y}) / f(\mathbf{y})$ . Given an observation vector  $\mathbf{Y} = \mathbf{y}$ , the decision rule is to allocate  $\mathbf{y}$  to population  $i$  if

$$\begin{aligned} P[J = i | \mathbf{Y} = \mathbf{y}] &> P[J = j | \mathbf{Y} = \mathbf{y}], \quad \text{for all } j \neq i, \text{ ie,} \\ \pi_i f_i(\mathbf{y}) &> \pi_j f_j(\mathbf{y}), \quad \text{for all } j \neq i. \end{aligned} \quad (6)$$

This is known as the Bayes' rule for classification, and it is the "best" rule as will be discussed below.

### 12.11.2 The Normal Case

We now discuss the case where the populations are  $N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) i = 1, \dots, k$ . If the populations have the same covariance matrix, that is,  $\boldsymbol{\Sigma}_1 = \dots = \boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}$ , then we are led to the linear discriminant rule. Otherwise, we have a quadratic discriminant rule.

**Case I:**  $\boldsymbol{\Sigma}_1 = \dots = \boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}$ .

The pdf of the  $j$ th population is

$$f_j(\mathbf{y}) = \left(1/\sqrt{2\pi}\right)^p (1/|\boldsymbol{\Sigma}|^{1/2}) \exp\left[-(1/2)(\mathbf{y} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}_j)\right].$$

According to the Bayes' rule given in Eq. (6), allocate  $\mathbf{y}$  to population  $i$  if, for all  $j \neq i$ ,

$$\begin{aligned} \pi_i f_i(\mathbf{y}) &> \pi_j f_j(\mathbf{y}), \text{ ie,} \\ -2 \log f_i(\mathbf{y}) - 2 \log \pi_i &< -2 \log f_j(\mathbf{y}) - 2 \log \pi_j, \text{ ie,} \\ (\mathbf{y} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}_i) - 2 \log \pi_i &< (\mathbf{y} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}_j) - 2 \log \pi_j. \end{aligned}$$

This is precisely the rule defined by the inequalities given in (5). If  $\pi_1, \dots, \pi_k$  are unknown, one often assumes that the prior is noninformative, that is,  $\pi_1 = \dots = \pi_k = 1/k$ . In such a case, we are led to the rule given in (4).

**Case II:**  $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k$  not the same.

The pdf of the  $j$ th population is

$$f_j(\mathbf{y}) = \left(1/\sqrt{2\pi}\right)^p (1/|\boldsymbol{\Sigma}_j|^{1/2}) \exp\left[-(1/2)(\mathbf{y} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{y} - \boldsymbol{\mu}_j)\right].$$

According to the Bayes' rule given in Eq. (6) allocate  $\mathbf{y}$  to population  $i$  if, for all  $j \neq i$ ,

$$\begin{aligned} \pi_i f_i(\mathbf{y}) &> \pi_j f_j(\mathbf{y}), \text{ ie,} \\ -2 \log f_i(\mathbf{y}) - 2 \log \pi_i &< -2 \log f_j(\mathbf{y}) - 2 \log \pi_j, \text{ ie,} \\ (\mathbf{y} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{y} - \boldsymbol{\mu}_i) + \log |\boldsymbol{\Sigma}_i| - 2 \log \pi_i &< (\mathbf{y} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{y} - \boldsymbol{\mu}_j) + \log |\boldsymbol{\Sigma}_j| - 2 \log \pi_j. \end{aligned}$$

This is called the quadratic discriminant rule. Note that unlike the case  $\Sigma_1 = \dots = \Sigma_k$ , the quadratic terms involving  $\mathbf{y}$  do not cancel out.

### 12.11.3 Sample Estimates

Suppose that we have  $n_i$  observations from population  $i$ ,  $i = 1, \dots, k$ . Then we can estimate the population means  $\mu_1, \dots, \mu_k$ . Let  $\mathbf{S}_i$  be the sample covariance matrix on the basis of  $n_i$  observations from population  $i$ .

#### *Linear Discriminant Rule*

In order to apply the linear discriminant rule, we need an estimate of  $\Sigma$  in addition to the estimates of  $\mu_1, \dots, \mu_k$ . Recall that an unbiased estimate of  $\Sigma$  is given by

$$\mathbf{S} = \mathbf{S}_{\text{pooled}} = \frac{1}{n-k} [(n_1 - 1)\mathbf{S}_1 + \dots + (n_k - 1)\mathbf{S}_k],$$

where  $n = n_1 + \dots + n_k$  is the total number of observation vectors.

If  $\pi_1, \dots, \pi_k$  are known (or if they are estimated), then the discriminant rule is: allocate  $\mathbf{y}$  to population  $i$  if, for all  $j \neq i$

$$\begin{aligned} (\mathbf{y} - \hat{\boldsymbol{\mu}}_i)^T \mathbf{S}^{-1} (\mathbf{y} - \hat{\boldsymbol{\mu}}_i) - 2 \log \pi_i &< (\mathbf{y} - \hat{\boldsymbol{\mu}}_j)^T \mathbf{S}^{-1} (\mathbf{y} - \hat{\boldsymbol{\mu}}_j) - 2 \log \pi_j, \text{ ie,} \\ -2\hat{\boldsymbol{\mu}}_i^T \mathbf{S}^{-1} \mathbf{y} + \hat{\boldsymbol{\mu}}_i^T \mathbf{S}^{-1} \hat{\boldsymbol{\mu}}_i - 2 \log \pi_i &< -2\hat{\boldsymbol{\mu}}_j^T \mathbf{S}^{-1} \mathbf{y} + \hat{\boldsymbol{\mu}}_j^T \mathbf{S}^{-1} \hat{\boldsymbol{\mu}}_j - 2 \log \pi_j. \end{aligned}$$

This rule is simplified when the prior is noninformative since the terms involving  $\{\log \pi_i\}$  cancel out.

*Remark 12.11.1.* In some cases it may be possible to estimate  $\pi_1, \dots, \pi_k$ . Suppose the observations in the sample are  $(J_t, \mathbf{Y}_t)$ ,  $t = 1, \dots, n$ , where  $J_t$  are iid. In such a case,  $n_i = \{\# \text{ of } J_t = i\}$  is random, and  $(n_1, \dots, n_k)$  is  $\text{Multinomial}(n; \pi_1, \dots, \pi_k)$ . We can then use the estimate  $\hat{\pi}_j = n_j/n$  (or  $\hat{\pi}_j = (n_j + 1/2)/(n + 1/2)$ ) of  $\pi_j$  in the classification rule.

#### *Quadratic Discriminant Rule*

The quadratic discriminant rule is: allocate  $\mathbf{y}$  to population  $i$  if, for all  $j \neq i$ ,

$$\begin{aligned} (\mathbf{y} - \hat{\boldsymbol{\mu}}_i)^T \mathbf{S}_i^{-1} (\mathbf{y} - \hat{\boldsymbol{\mu}}_i) + \log |\mathbf{S}_i| - 2 \log \pi_i \\ < (\mathbf{y} - \hat{\boldsymbol{\mu}}_j)^T \mathbf{S}_j^{-1} (\mathbf{y} - \hat{\boldsymbol{\mu}}_j) + \log |\mathbf{S}_j| - 2 \log \pi_j. \end{aligned}$$

One big drawback of the quadratic discriminant rule is that we need to estimate  $\Sigma_1, \dots, \Sigma_k$  which is equivalent to estimating  $kp(p+1)/2$  parameters of the covariance matrices. This can lead to inefficiencies especially if  $p(p+1)/2$  is not small in comparison to  $\min\{n_1, \dots, n_k\}$ . A plausible remedy is to estimate  $\Sigma_i$  by shrinking  $\mathbf{S}_i$  toward  $\mathbf{S}$ , the pooled estimate constructed in the linear discriminant rule (ie, estimate  $\Sigma_i$  by  $\hat{\Sigma}_i = (1-\alpha_i)\mathbf{S}_i + \alpha_i\mathbf{S}$ ,  $0 \leq \alpha_i \leq 1$ ). One can then carry out a quadratic discriminant rule using  $\{\hat{\Sigma}_i\}$  instead of

$\{S_i\}$ . The constants  $\{\alpha_i\}$  need to be estimated from the data, and methods for doing this include:

- (i) minimizing an AIC-type criterion, and
- (ii) minimizing a cross-validation type criterion.

#### 12.11.4 Probability of Misclassification

Let  $\mathbf{Y}$  be a  $p$ -dim random vector which comes from one of the  $k$  populations with pdf's  $f_1, \dots, f_k$ . As before assume that  $\pi_1, \dots, \pi_k$  are the prior probabilities which may or may not be noninformative. If  $d(\mathbf{Y})$  is a function from  $\mathbb{R}^p$  to the set  $\{1, \dots, k\}$ , then it is called a classifier or a classification function. If  $d(\mathbf{Y}) = i$ , then the classifier  $d$  allocates  $\mathbf{Y}$  to population  $i$ . Now if  $\mathbf{Y}$  is actually from population  $i$ , but the classifier allocates it to population  $j$  (ie,  $d(\mathbf{Y}) = j, j \neq i$ ), then there is a misclassification. Given that  $\mathbf{Y}$  is from population  $i$ , the probability of misclassification is

$$P[d(\mathbf{Y}) \neq i | J = i] = \sum_{j \neq i} P[d(\mathbf{Y}) = j | J = i].$$

So the total probability of misclassification for the classifier is

$$\begin{aligned} P[d(\mathbf{Y}) \neq J] &= \sum_{i=1}^k P[d(\mathbf{Y}) \neq i | J = i] P[J = i] \\ &= \sum_{i=1}^k P[d(\mathbf{Y}) \neq i | J = i] \pi_i. \end{aligned}$$

The following result states that the classification rule defined by the inequalities given in (6) has the smallest total probability of misclassification among all classifiers.

**Theorem 12.11.1.** *Let  $d^*$  be the Bayes' rule for classification defined by the inequalities given in (6) that is,*

$$d^*(\mathbf{y}) = i \text{ if } \pi_i f_i(\mathbf{y}) > \pi_j f_j(\mathbf{y}) \text{ for all } j \neq i, i = 1, \dots, k.$$

*If  $\mathcal{D}$  is the set of all classifiers from  $\mathbb{R}^p$  to  $\{1, \dots, k\}$ , then  $\min_{d \in \mathcal{D}} P[d(\mathbf{Y}) \neq J] = P[d^*(\mathbf{Y}) \neq J]$ .*

*Remark 12.11.2.* Note that when  $\pi_i f_i(\mathbf{y}) = \pi_j f_j(\mathbf{y})$  for some  $i \neq j$ , there may be an ambiguity in how to classify  $\mathbf{y}$ . If  $\mathbf{Y}$  has a continuous distribution, then  $P[\pi_i f_i(\mathbf{Y}) = \pi_j f_j(\mathbf{Y})] = 0$  for any  $j \neq i$ . For this reason, the definition of  $d^*$  as given above is adequate for continuous distributions.

*Proof of Theorem 12.11.1.* Since  $P[d(\mathbf{Y}) \neq J] = 1 - P[d(\mathbf{Y}) = J]$ , it is enough to show that for any classification rule  $d$ ,

$$P[d(\mathbf{Y}) = J] \leq P[d^*(\mathbf{Y}) = J].$$

Note that  $P[d(\mathbf{Y}) = J] = E[P\{d(\mathbf{Y}) = J | \mathbf{Y}\}]$ . We will show that for classification rule  $d$

$$P[d(\mathbf{Y}) = J | \mathbf{Y}] \leq P[d^*(\mathbf{Y}) = J | \mathbf{Y}].$$

In the calculation of the conditional probability  $P[d(\mathbf{Y}) = J|\mathbf{Y}]$ , we may assume that  $d(\mathbf{Y})$  is fixed since  $\mathbf{Y}$  is fixed. Hence

$$P[d(\mathbf{Y}) = J|\mathbf{Y}] = \sum_{i=1}^k I(d(\mathbf{Y}) = i)P[J = i|\mathbf{Y}],$$

where  $I$  is the indicator function such that  $I(u = v) = 1$  if  $u = v$  and  $= 0$  otherwise. Assume that the maximum of  $P[J = i|\mathbf{Y}]$ , over  $i = 1, \dots, k$ , is attained at  $i^*$ . Hence by definition,  $d^*(\mathbf{y}) = i^*$ . So  $I(d^*(\mathbf{Y}) = i)$  is equal to 0 or 1 depending on whether  $i \neq i^*$  or  $i = i^*$ . Hence

$$\begin{aligned} P[d(\mathbf{Y}) = J|\mathbf{Y}] &= \sum_{i=1}^k I(d(\mathbf{Y}) = i)P[J = i|\mathbf{Y}] \\ &\leq \sum_{i=1}^k I(d(\mathbf{Y}) = i)P[J = i^*|\mathbf{Y}] = P[J = i^*|\mathbf{Y}] \\ &= \sum_{i=1}^k I(d^*(\mathbf{Y}) = i)P[J = i|\mathbf{Y}] = P[d^*(\mathbf{Y}) = J]. \end{aligned}$$

This concludes the proof of this result. □

### 12.11.5 Classification: Fisher's Method

This method does not require the normality assumption and is flexible enough to provide nonlinear classification rules. Suppose that we have  $n_i$  observations vectors  $\{\mathbf{Y}_{ij}: j = 1, \dots, n_i\}$  from population  $i$ ,  $i = 1, \dots, k$ . If  $\mathbf{e}$  is in  $\mathbb{R}^p$ , then we have a one-factor ANOVA model

$$\mathbf{e}^T \mathbf{Y}_{ij} = \mathbf{e}^T \boldsymbol{\mu}_i + \mathbf{e}^T \boldsymbol{\epsilon}_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

and  $R^2(\mathbf{e})$ , the coefficient of determination (given in Eq. (1)), is

$$R^2(\mathbf{e}) = \frac{\mathbf{e}^T \mathbf{B} \mathbf{e}}{\mathbf{e}^T \mathbf{T} \mathbf{e}} = \frac{\mathbf{e}^T \mathbf{B} \mathbf{e}}{\mathbf{e}^T \mathbf{B} \mathbf{e} + \mathbf{e}^T \mathbf{W} \mathbf{e}} = \frac{(\mathbf{e}^T \mathbf{B} \mathbf{e}) / (\mathbf{e}^T \mathbf{W} \mathbf{e})}{(\mathbf{e}^T \mathbf{B} \mathbf{e}) / (\mathbf{e}^T \mathbf{W} \mathbf{e}) + 1},$$

where  $\mathbf{B}$ ,  $\mathbf{W}$ , and  $\mathbf{T}$  are the between group, within group, and total SSP matrices, respectively. Maximizing  $R^2(\mathbf{e})$  with respect to  $\mathbf{e}$  leads to a generalized eigenvalue problem (Section B.4). Now maximizing  $R^2(\mathbf{e})$  is equivalent to maximizing the ratio  $\mathbf{e}^T \mathbf{B} \mathbf{e} / \mathbf{e}^T \mathbf{W} \mathbf{e}$  which in turn is equivalent to maximizing  $\mathbf{e}^T \mathbf{B} \mathbf{e} / \mathbf{e}^T \mathbf{S} \mathbf{e}$ , where  $\mathbf{S} = (n - k)^{-1} \mathbf{W}$  is the pooled covariance matrix ( $n = n_1 + \dots + n_k$ ).

Let  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  be the eigenvalues of  $\mathbf{B}$  with respect to  $\mathbf{S}$  (ie, these are the eigenvalues of  $\mathbf{S}^{-1/2} \mathbf{B} \mathbf{S}^{-1/2}$ ). Let  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$  be the corresponding orthonormal eigenvectors of  $\mathbf{S}^{-1/2} \mathbf{B} \mathbf{S}^{-1/2}$ . If  $\hat{\mathbf{e}}_j = \mathbf{S}^{-1/2} \hat{\mathbf{u}}_j$ ,  $j = 1, 2, \dots$ , then  $\hat{\mathbf{e}}_j^T \mathbf{S} \hat{\mathbf{e}}_j = 1$  for all  $j$ , and  $\hat{\mathbf{e}}_i^T \mathbf{S} \hat{\mathbf{e}}_j = 0$  whenever  $i \neq j$ . Since  $s = \text{rank}(\mathbf{B}) = \min(k - 1, p)$ ,  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_s > 0$  and  $\hat{\lambda}_{s+1} = \dots = \hat{\lambda}_p = 0$ . We will call  $Z_1 = \hat{\mathbf{e}}_1^T \mathbf{Y}$  the first discriminant,  $Z_2 = \hat{\mathbf{e}}_2^T \mathbf{Y}$  the second discriminant, and so on. The

vector of discriminants  $\mathbf{Z}$  is given by  $(\hat{\mathbf{e}}_1^T \mathbf{Y}, \dots, \hat{\mathbf{e}}_s^T \mathbf{Y})^T$ . If we write  $\mathbf{Z}_{ij} = (\hat{\mathbf{e}}_1^T \mathbf{Y}_{ij}, \dots, \hat{\mathbf{e}}_s^T \mathbf{Y}_{ij})^T$ , then we can write an approximate MANOVA model

$$\mathbf{Z}_{ij} = \tilde{\boldsymbol{\theta}}_i + \delta_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k,$$

where  $\tilde{\boldsymbol{\theta}}_i = (\hat{\mathbf{e}}_1^T \boldsymbol{\mu}_i, \dots, \hat{\mathbf{e}}_s^T \boldsymbol{\mu}_i)^T$ , and for each  $i$ , sample means and sample covariance matrix of  $\{\mathbf{Z}_{ij}\}$  are  $\{\tilde{\boldsymbol{\theta}}_i\}$  and  $\mathbf{I}$ , respectively.

### Method of Classification

Suppose  $\mathbf{y}$  is from one of the  $k$  populations and we need to classify it. Let  $\mathbf{z} = (\hat{\mathbf{e}}_1^T \mathbf{y}, \dots, \hat{\mathbf{e}}_s^T \mathbf{y})^T$  be the vector of discriminant scores of  $\mathbf{y}$ . Let  $\hat{\boldsymbol{\theta}}_i = (\hat{\mathbf{e}}_1^T \hat{\boldsymbol{\mu}}_i, \dots, \hat{\mathbf{e}}_s^T \hat{\boldsymbol{\mu}}_i)^T$ . Allocate  $\mathbf{y}$  to population  $i$  if

$$\|\mathbf{z} - \hat{\boldsymbol{\theta}}_i\|^2 < \|\mathbf{z} - \hat{\boldsymbol{\theta}}_j\|^2 \quad \text{for all } j \neq i.$$

If it is decided to use  $q$  discriminants ( $q < s$ ), then we create a vector  $\mathbf{z} = (\hat{\mathbf{e}}_1^T \mathbf{y}, \dots, \hat{\mathbf{e}}_q^T \mathbf{y})^T$  using the first  $q$  discriminants and let  $\hat{\boldsymbol{\theta}}_i = (\hat{\mathbf{e}}_1^T \hat{\boldsymbol{\mu}}_i, \dots, \hat{\mathbf{e}}_q^T \hat{\boldsymbol{\mu}}_i)^T$  be the corresponding sample means. Then the rule is: allocate  $\mathbf{y}$  to population  $i$  if

$$\|\mathbf{z} - \hat{\boldsymbol{\theta}}_i\|^2 < \|\mathbf{z} - \hat{\boldsymbol{\theta}}_j\|^2 \quad \text{for all } j \neq i.$$

In practice, the decision to use the first  $q$  discriminants is usually based on how close the ratio  $\hat{\pi}_q = (\hat{\lambda}_1 + \dots + \hat{\lambda}_q) / (\hat{\lambda}_1 + \dots + \hat{\lambda}_s)$  is to 1, and this approach makes intuitive sense since  $\hat{\pi}_q$  is an estimate of the proportion of variability in  $\mathbf{Y}$  explained by the first  $q$  discriminants.

### Connection to Linear Discriminant Rule

The linear discriminant rule is not necessarily the same as Fisher's. These two rules are the same if all the discriminants are used (ie, number of discriminants is equal to  $s = \text{rank}(\mathbf{B})$ ) in Fisher's method and the prior is noninformative.

**Lemma 12.11.1.** *Fisher's classification rule with  $s = \text{rank}(\mathbf{B})$  discriminants is equivalent to the linear discriminant rule with a noninformative prior (ie,  $\pi_1 = \dots = \pi_k = 1/k$ ).*

*Proof of Lemma 12.11.1.* Recall that  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  are the eigenvalues of  $\mathbf{S}^{-1/2} \mathbf{B} \mathbf{S}^{-1/2}$  with the corresponding orthonormal eigenvectors  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$ . Set  $\hat{\mathbf{e}}_l = \mathbf{S}^{-1/2} \hat{\mathbf{u}}_l$ ,  $l = 1, \dots, p$ . Since  $\hat{\lambda}_l = 0$  for  $l \geq s+1$ , we have for any  $l \geq s+1$

$$0 = \hat{\mathbf{u}}_l^T \mathbf{S}^{-1/2} \mathbf{B} \mathbf{S}^{-1/2} \hat{\mathbf{u}}_l = \hat{\mathbf{e}}_l^T \mathbf{B} \hat{\mathbf{e}}_l = \sum_{1 \leq i \leq k} n_i (\hat{\mathbf{e}}_l^T \bar{\mathbf{Y}}_i - \hat{\mathbf{e}}_l^T \bar{\mathbf{Y}}..)^2.$$

Hence,  $\hat{\mathbf{e}}_l^T \bar{\mathbf{Y}}_i = \hat{\mathbf{e}}_l^T \bar{\mathbf{Y}}..$  whenever  $l \geq s+1$ . For  $1 \leq i \leq k$ , and any  $\mathbf{y}$  in  $\mathbb{R}^p$ ,

$$\begin{aligned} & (\mathbf{y} - \bar{\mathbf{Y}}_i)^T \mathbf{S}^{-1} (\mathbf{y} - \bar{\mathbf{Y}}_i) \\ &= \|\mathbf{S}^{-1/2} (\mathbf{y} - \bar{\mathbf{Y}}_i)\|^2 \\ &= \sum_{l=1}^p (\hat{\mathbf{u}}_l^T \mathbf{S}^{-1/2} (\mathbf{y} - \bar{\mathbf{Y}}_i))^2 \quad (\{\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_p\} \text{ is an orthonormal basis of } \mathbb{R}^p) \end{aligned}$$

$$\begin{aligned}
&= \sum_{l=1}^p \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2 = \sum_{1 \leq l \leq s} \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2 + \sum_{s+1 \leq l \leq p} \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2 \\
&= \sum_{l=1}^s \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2 + \sum_{l=s+1}^p \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2 \\
&= \|\mathbf{z} - \hat{\boldsymbol{\theta}}_i\|^2 + \sum_{l=s+1}^p \{\hat{\mathbf{e}}_l^T (\mathbf{y} - \bar{\mathbf{Y}}_{i.})\}^2,
\end{aligned}$$

where  $\mathbf{z} = (\hat{\mathbf{e}}_1^T \mathbf{y}, \dots, \hat{\mathbf{e}}_s^T \mathbf{y})^T$  and  $\hat{\boldsymbol{\theta}}_i = (\hat{\mathbf{e}}_1^T \hat{\boldsymbol{\mu}}_i, \dots, \hat{\mathbf{e}}_s^T \hat{\boldsymbol{\mu}}_i)^T$ . Note that the second sum in the last expression does not depend on  $i$ . So comparing  $(\mathbf{y} - \bar{\mathbf{Y}}_{i.})^T \mathbf{S}^{-1} (\mathbf{y} - \bar{\mathbf{Y}}_{i.})$  to  $(\mathbf{y} - \bar{\mathbf{Y}}_{i'.})^T \mathbf{S}^{-1} (\mathbf{y} - \bar{\mathbf{Y}}_{i'})$ ,  $i \neq i'$ , is equivalent to comparing  $\|\mathbf{z} - \hat{\boldsymbol{\theta}}_i\|^2$  to  $\|\mathbf{z} - \hat{\boldsymbol{\theta}}_{i'}\|^2$ .

This argument shows that, for the noninformative prior, the linear discriminant rule is equivalent to Fisher's rule if the number of discriminants used by Fisher's method is equal to  $\text{rank}(\mathbf{B})$ .

□

## 12.12 Canonical Correlation Analysis

Canonical correlation analysis is a descriptive method that seeks to obtain measures of association between two sets of multivariate observations. Let  $\mathbf{X}$  be  $q \times 1$  and  $\mathbf{Y}$  be  $p \times 1$  random vectors with means  $\boldsymbol{\mu}_X$  and  $\boldsymbol{\mu}_Y$ , respectively. Assume that the covariance matrix of  $\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$  is  $\begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$ ,

$$\boldsymbol{\Sigma}_{11} = \text{Cov}[\mathbf{X}], \boldsymbol{\Sigma}_{22} = \text{Cov}[\mathbf{Y}], \text{ and } \boldsymbol{\Sigma}_{12} = \text{Cov}[\mathbf{X}, \mathbf{Y}] = \boldsymbol{\Sigma}_{21}^T.$$

Assume that  $\boldsymbol{\Sigma}_{11}$  and  $\boldsymbol{\Sigma}_{22}$  are nonsingular.

The goal is to find linear functions  $\mathbf{a}^T \mathbf{X}$  and  $\mathbf{b}^T \mathbf{Y}$ ,  $\mathbf{a} \in \mathbb{R}^q$  and  $\mathbf{b} \in \mathbb{R}^p$ , which maximize the correlation between  $\mathbf{a}^T \mathbf{X}$  and  $\mathbf{b}^T \mathbf{Y}$ . Clearly, we can rescale  $\mathbf{a}$  and  $\mathbf{b}$  so that  $\text{Var}[\mathbf{a}^T \mathbf{X}] = \text{Var}[\mathbf{b}^T \mathbf{Y}] = 1$  (ie,  $\mathbf{a}^T \boldsymbol{\Sigma}_{11} \mathbf{a} = \mathbf{b}^T \boldsymbol{\Sigma}_{22} \mathbf{b} = 1$ ). In such a case,  $\text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] = \mathbf{a}^T \boldsymbol{\Sigma}_{12} \mathbf{b}$ .

If  $\text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] \leq 0$ , then the correlation between  $-\mathbf{a}^T \mathbf{X}$  and  $\mathbf{b}^T \mathbf{Y}$  is nonnegative. This shows that maximizing the absolute value of  $\text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}]$  is equivalent to maximizing  $\text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}]$ .

Let  $\lambda_1 \geq \lambda_2 \geq \dots$  be the eigenvalues of the matrix  $\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1}$  with  $\mathbf{u}_1, \mathbf{u}_2, \dots$  as the corresponding orthonormal eigenvectors. Let

$$\rho_i = \lambda_i^{1/2}, \quad \mathbf{a}_i = \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{u}_i, \text{ and } \mathbf{b}_i = \rho_i^{-1} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_i, \tag{7}$$

$i = 1, \dots, k = \min(p, q)$ . Then  $\rho_1 \geq \rho_2 \geq \dots$  are called the canonical correlations,  $\mathbf{a}_1^T \mathbf{X}, \dots$  are called the canonical variates of  $\mathbf{X}$ , and  $\mathbf{b}_1^T \mathbf{Y}, \mathbf{b}_2^T \mathbf{Y}, \dots$  are called the canonical variates of  $\mathbf{Y}$ . Detailed arguments are given in [Section 12.12.4](#).

The following result summarizes the key ideas.

**Theorem 12.12.1.** Let  $\{\rho_i\}$ ,  $\{\mathbf{a}_i\}$ , and  $\{\mathbf{b}_i\}$  be as given in Eq. (7). Then the following hold:

- (a)  $\text{Var}[\mathbf{a}_i^T \mathbf{X}] = \text{Var}[\mathbf{b}_i^T \mathbf{Y}] = 1, i = 1, \dots, k.$
- (b) When  $i \neq j$ ,  $\text{Cov}[\mathbf{a}_i^T \mathbf{X}, \mathbf{a}_j^T \mathbf{X}] = 0$ ,  $\text{Cov}[\mathbf{b}_i^T \mathbf{Y}, \mathbf{b}_j^T \mathbf{Y}] = 0$ , and  $\text{Cov}[\mathbf{a}_i^T \mathbf{X}, \mathbf{b}_j^T \mathbf{Y}] = 0$ .
- (c)  $\text{Corr}[\mathbf{a}_i^T \mathbf{X}, \mathbf{b}_i^T \mathbf{Y}] = \rho_i, i = 1, \dots, k.$

Proof of this result is given in [Section 12.12.4](#). We now write down two more results on canonical correlations.

**Theorem 12.12.2.** Let  $\{\rho_i\}$ ,  $\{\mathbf{a}_i\}$ , and  $\{\mathbf{b}_i\}$  be as in [Theorem 12.12.1](#). Then:

- (a)  $\max_{\mathbf{a}^T \Sigma_{11} \mathbf{a}=1, \mathbf{b}^T \Sigma_{22} \mathbf{b}=1} \text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] = \text{Corr}[\mathbf{a}_1^T \mathbf{X}, \mathbf{b}_1^T \mathbf{Y}] = \rho_1.$
- (b) Under the constraints  $\mathbf{a}^T \Sigma_{11} \mathbf{a}_i, i = 1, \dots, r - 1$ ,  
 $\max_{\mathbf{a}^T \Sigma_{11} \mathbf{a}=1, \mathbf{b}^T \Sigma_{22} \mathbf{b}=1} \text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] = \text{Corr}[\mathbf{a}_r^T \mathbf{X}, \mathbf{b}_r^T \mathbf{Y}] = \rho_r, r = 2, \dots, k.$

Proof of [Theorem 12.12.2](#) is not given since it is basically a restatement of [Theorem B.4.1](#) given in [Appendix B](#). The following result on invariance of canonical correlations under nonsingular linear transformations is left as an exercise.

**Lemma 12.12.1** (Invariance). If  $\tilde{\mathbf{X}} = \mathbf{c} + \mathbf{U}\mathbf{X}$  and  $\tilde{\mathbf{Y}} = \mathbf{d} + \mathbf{V}\mathbf{Y}$ , where  $\mathbf{U}$  is  $q \times q$  and nonsingular,  $\mathbf{V}$  is  $p \times p$  and nonsingular, and  $\mathbf{c}$ ,  $\mathbf{d}$ ,  $\mathbf{U}$ , and  $\mathbf{V}$  are nonrandom, then the canonical correlations between  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$  are the same as those between  $\mathbf{X}$  and  $\mathbf{Y}$ . Moreover, if  $\{\mathbf{a}_i\}$  and  $\{\mathbf{b}_i\}$  are the canonical vectors of  $\mathbf{X}$  and  $\mathbf{Y}$ , then  $\tilde{\mathbf{a}}_i = \{\mathbf{U}^{-1}\}^T \mathbf{a}_i$  and  $\tilde{\mathbf{b}}_i = \{\mathbf{V}^{-1}\}^T \mathbf{b}_i$  are the canonical vectors of  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$ , respectively,  $i = 1, \dots, k$ .

### 12.12.1 Sample Estimates

Let  $(\mathbf{X}_t, \mathbf{Y}_t)$ ,  $t = 1, \dots, n$ , be the  $n$  pairs of vector observations, and let  $\mathbf{S}_{11}$  and  $\mathbf{S}_{22}$  be the sample covariance matrices of  $\{\mathbf{X}_t\}$  and  $\{\mathbf{Y}_t\}$ , respectively. Set  $\mathbf{S}_{12} = \frac{1}{n-1} \sum (\mathbf{X}_t - \bar{\mathbf{X}})(\mathbf{Y}_t - \bar{\mathbf{Y}})^T$  and  $\mathbf{S}_{21} = \mathbf{S}_{12}^T$ .

Let  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  be the eigenvalues of the matrix  $\mathbf{S}_{11}^{-1/2} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1}$  with  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$  as the corresponding orthonormal eigenvectors. Let

$$\hat{\rho}_i = \hat{\lambda}_i^{1/2}, \quad \hat{\mathbf{a}}_i = \mathbf{S}_{11}^{-1/2} \hat{\mathbf{u}}_i, \text{ and } \hat{\mathbf{b}}_i = \hat{\rho}_i^{-1} \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \hat{\mathbf{a}}_i,$$

$i = 1, \dots, k = \min(p, q)$ . Then

$\hat{\rho}_1 \geq \hat{\rho}_2 \geq \dots$  are the estimated canonical correlations,  
 $\hat{\mathbf{a}}_1^T \mathbf{X}, \hat{\mathbf{a}}_2^T \mathbf{X}, \dots$  are the estimated canonical variates of  $\mathbf{X}$ , and  
 $\hat{\mathbf{b}}_1^T \mathbf{Y}, \hat{\mathbf{b}}_2^T \mathbf{Y}, \dots$  are the estimated canonical variates of  $\mathbf{Y}$ .

### 12.12.2 Test for $\Sigma_{12} = \mathbf{0}$

Are  $\mathbf{X}$  and  $\mathbf{Y}$  uncorrelated? The likelihood ratio statistic for testing  $H_0: \Sigma_{12} = \mathbf{0}$  against  $H_1: \Sigma_{12} \neq \mathbf{0}$  is

$$\lambda = |\mathbf{I} - \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12}|^{n/2} = \{(1 - \hat{\rho}_1^2) \cdots (1 - \hat{\rho}_k^2)\}^{n/2} = \Lambda^{n/2},$$

where  $\Lambda = (1 - \hat{\rho}_1^2) \cdots (1 - \hat{\rho}_k^2)$  is Wilks' lambda. Under  $H_0$ ,

$$-[n - (p + q + 3)/2] \log \Lambda \xrightarrow{\mathcal{D}} \chi_{pq}^2,$$

as  $n \rightarrow \infty$ .

If we want to test that only  $m < k$  of the population canonical correlations are non-negative, then the test is based on the statistic  $\Lambda = (1 - \hat{\rho}_{m+1}^2) \cdots (1 - \hat{\rho}_k^2)$ . Then under  $H_0$ :  $\rho_{m+1} = \cdots = 0$ ,

$$-[n - (p + q + 3)/2] \log \Lambda \xrightarrow{\mathcal{D}} \chi_{(p-m)(q-m)}^2,$$

as  $n \rightarrow \infty$ .

**Example 12.12.1.** Scores of  $n = 88$  students in five subjects are given in Mardia et al. [61]. The subjects are Mechanics (C), Vectors (C), Algebra (O), Analysis (O), and Statistics (O).

Here “C” and “O” stand for closed- and open-book examinations. The goal is to find canonical correlations between open- and closed-book scores. Let  $\mathbf{X}$  be the vector closed-book marks and  $\mathbf{Y}$  be the vector of open-book marks. The sample covariance matrix is

$$\mathbf{S} = \begin{pmatrix} 302.3 & 125.8 & 100.4 & 105.1 & 116.1 \\ & 170.9 & 84.2 & 93.6 & 97.9 \\ & & 111.6 & 110.8 & 120.5 \\ & & & 217.9 & 153.8 \\ & & & & 294.4 \end{pmatrix}.$$

Here  $q = 2$  and  $p = 3$ , and so  $k = \min(q, p) = 2$ .

The canonical correlations are  $\hat{\rho}_1 = 0.6630$  and  $\hat{\rho}_2 = 0.0412$ . The first canonical scores are

$$\hat{\mathbf{a}}_1^T \mathbf{X} = 0.0260X_1 + 0.0518X_2, \text{ and } \hat{\mathbf{b}}_1^T \mathbf{Y} = 0.0824Y_1 + 0.0081Y_2 + 0.0035Y_3.$$

### 12.12.3 Cross-Classified Data and Canonical Correlation

Consider the  $5 \times 5$  contingency table given below on the occupational status of  $n = 3497$  father-son pairs. Data sets of this type have been analyzed by many authors in order to investigate issues of intergenerational mobility. Note that occupational status is a qualitative variable and its numbering of 1 through 5 is purely descriptive. Can we assign numerical values for father's and son's status so that they can be treated as “quantitative” variables? This issue will be addressed here using canonical correlation analysis.

Suppose that we have an  $r \times c$  contingency table with cell counts  $\{n_{ij}\}$  and the total is  $n = \sum \sum n_{ij}$ . For  $i = 1, \dots, r$ , and  $j = 1, \dots, c$ , create the following indicator variables,  $t = 1, \dots, n$ ,

$$X_{ti} = \begin{cases} 1 & \text{if individual } t \text{ belongs to the } i\text{th row category} \\ 0 & \text{otherwise} \end{cases}$$

$$Y_{tj} = \begin{cases} 1 & \text{if individual } t \text{ belongs to the } j\text{th column category} \\ 0 & \text{otherwise.} \end{cases}$$

Thus the observations are  $(X_t, Y_t)$ ,  $t = 1, \dots, n$ , where  $X_t = (X_{t1}, \dots, X_{tr})^T$  and  $Y_t = (Y_{t1}, \dots, Y_{tc})^T$ . Suppose the numerical values for father's status and son's status are  $a_1, \dots, a_r$  and  $b_1, \dots, b_c$ , respectively (here  $r = c = 5$ ). We wish to find  $\{a_j\}$  and  $\{b_j\}$  so that the correlation between father's and son's status is maximized and this is done by employing canonical correlation analysis. Since  $X$  has only one nonzero component,  $\mathbf{a}^T X$  takes one of the  $r$  values  $a_1, \dots, a_r$ . Similarly,  $\mathbf{b}^T Y$  takes one of the  $c$  values  $b_1, \dots, b_c$ . Thus a father-son pair has a bivariate score  $(a_i, b_j)$ , if the father is in the  $i$ th row category and the son is in the  $j$ th column category. We can now find  $\mathbf{a} = (a_1, \dots, a_r)^T$  and  $\mathbf{b} = (b_1, \dots, b_c)^T$  to maximize  $\text{Corr}[\mathbf{a}^T X, \mathbf{b}^T Y]$ .

Social mobility data:  $n = 3497$

Father's Status	Subject's Status					Total	Percent
	1	2	3	4	5		
1	50	45	8	18	8	129	3.7
2	28	174	84	154	55	495	14.2
3	11	78	110	223	96	516	14.8
4	14	150	185	714	447	1510	43.2
5	0	42	72	320	411	845	24.2
Total	103	489	459	1429	1017	3497	
Percent	2.9	13.9	13.1	40.9	29.1		

Classes: 1 = professional, 2 = intermediate, 3 = skilled, 4 = semiskilled, 5 = unskilled.

The first canonical correlation is  $\hat{\rho}_1 = 0.504$ . The canonical scores are given below

Father's Status	1	2	3	4	5
	0	3.15	4.12	4.55	4.96

Son's Status	1	2	3	4	5
	0	3.34	4.49	4.87	5.26

Scores seem to be increasing for both the father and the son. Father's scores seem to be correlated to son's scores. Social Classes 1 and 2 seem to be more distinct from one another than other adjacent social classes, both for the son's and the father's.

### 12.12.4 Technical Notes

#### *Derivation of Canonical Correlations*

The following lemma is useful and it is similar to a result in [Section B.1](#).

**Lemma 12.12.2.** *Let  $\mathbf{w}$  be in  $\mathbb{R}^m$ . Then  $\max_{\|\mathbf{a}\|=1} \mathbf{w}^T \mathbf{a} = \|\mathbf{w}\|$ , and the maximum is attained at  $\mathbf{a} = \mathbf{w}/\|\mathbf{w}\|$ .*

Now let us find the canonical correlations and the corresponding canonical vectors. Since  $\mathbf{a}^T \Sigma_{11} \mathbf{a} = \mathbf{b}^T \Sigma_{22} \mathbf{b} = 1$ , we have

$$\text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] = \text{Cov}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] = \mathbf{a}^T \Sigma_{12} \mathbf{b}.$$

Making a change of vector  $\mathbf{v} = \Sigma_{22}^{1/2} \mathbf{b}$ , we have

$$\begin{aligned} \max_{\mathbf{b}^T \Sigma_{22} \mathbf{b}=1} \text{Corr}[\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}] &= \max_{\|\mathbf{v}\|=1} \mathbf{a}^T \Sigma_{12} \Sigma_{22}^{-1/2} \mathbf{v} \\ &= \|\Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}\| = (\mathbf{a}^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a})^{1/2}, \end{aligned}$$

and this maximum occurs at  $\mathbf{v} = \Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}/\|\Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}\|$ , that is, at  $\mathbf{b} = \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}/\|\Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}\|$ .

Now making a change of vector  $\mathbf{u} = \Sigma_{11}^{1/2} \mathbf{a}$ , we get

$$\max_{\mathbf{a}^T \Sigma_{11} \mathbf{a}=1} \mathbf{a}^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a} = \max_{\|\mathbf{u}\|=1} \mathbf{u}^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{u}.$$

Let  $\lambda_1 \geq \lambda_2 \geq \dots$  be the eigenvalues of  $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$  with  $\mathbf{u}_1, \mathbf{u}_2, \dots$  as the corresponding orthonormal eigenvectors. So

$$\rho_i = \lambda_i^{1/2}, \quad \mathbf{a}_i = \Sigma_{11}^{-1/2} \mathbf{u}_i, \text{ and } \mathbf{b}_i = \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}_i / \|\Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}_i\|.$$

Note that

$$\|\Sigma_{22}^{-1/2} \Sigma_{21} \mathbf{a}_i\|^2 = \mathbf{a}_i^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}_i = \mathbf{u}_i^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{u}_i = \rho_i^2.$$

Thus we have  $\mathbf{b}_i = \rho_i^{-1} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}_i$ .

*Proof of Theorem 12.12.1.*

**(a)** It is easy to see that  $\text{Var}[\mathbf{a}_i^T \mathbf{X}] = \mathbf{a}_i^T \Sigma_{11} \mathbf{a}_i = \mathbf{u}_i^T \mathbf{u}_i = 1$ .

Note that

$$\begin{aligned} \text{Var}[\mathbf{b}_i^T \mathbf{Y}] &= \mathbf{b}_i^T \Sigma_{22} \mathbf{b}_i = \rho_i^{-2} \mathbf{a}_i^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12} \mathbf{a}_i \\ &= \rho_i^{-2} \mathbf{u}_i^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{u}_i = \rho_i^{-2} \rho_i^2 = 1. \end{aligned}$$

**(b)** It is fairly easy to see that  $\mathbf{a}_i^T \mathbf{X}$  and  $\mathbf{a}_j^T \mathbf{X}$  are uncorrelated for  $i \neq j$  since

$$\text{Cov}[\mathbf{a}_i^T \mathbf{X}, \mathbf{a}_j^T \mathbf{X}] = \mathbf{a}_i^T \Sigma_{11} \mathbf{a}_j = \mathbf{u}_i^T \mathbf{u}_j = 0.$$

Note that when  $i \neq j$ , we have

$$\begin{aligned} \text{Cov}[\mathbf{b}_i^T \mathbf{Y}, \mathbf{b}_j^T \mathbf{Y}] &= \mathbf{b}_i^T \Sigma_{22} \mathbf{b}_j = \rho_i^{-1} \rho_j^{-1} \mathbf{a}_i^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}_j \\ &= \rho_i^{-1} \rho_j^{-1} \mathbf{u}_i^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{u}_j = 0. \end{aligned}$$

The last step follows from the fact that the eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots$  are orthonormal.  
Similarly

$$\text{Cov}[\mathbf{a}_i^T \mathbf{X}, \mathbf{b}_j^T \mathbf{Y}] = \mathbf{a}_i^T \boldsymbol{\Sigma}_{12} \mathbf{b}_j = \rho_j^{-1} \mathbf{a}_i^T \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_j = 0.$$

The last step follows from the intermediate step in the proof for  $\text{Cov}[\mathbf{b}_i^T \mathbf{Y}, \mathbf{b}_j^T \mathbf{Y}] = 0$ .

(c) Note that

$$\begin{aligned}\text{Corr}[\mathbf{a}_i^T \mathbf{X}, \mathbf{b}_i^T \mathbf{Y}] &= \text{Cov}[\mathbf{a}_i^T \mathbf{X}, \mathbf{b}_i^T \mathbf{Y}] \\ &= \mathbf{a}_i^T \boldsymbol{\Sigma}_{12} \mathbf{b}_i = \rho_i^{-1} \mathbf{a}_i^T \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \mathbf{a}_i \\ &= \rho_i^{-1} \mathbf{u}_i^T \boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{u}_i = \rho_i^{-1} \rho_i^2 = \rho_i.\end{aligned}$$

□

## Exercises

- 12.1.** For this question, you may use the following facts: if  $X \sim \chi_m^2$ , then  $E[X] = m$  and  $E[1/X] = 1/(m - 2)$ ,  $m > 2$
- (a) Show that  $E[F_{u,v}] = \frac{v}{v-2}$  where  $F_{u,v}$  has an  $F$ -distribution with  $df = (u, v)$ ,  $v > 2$ . [Hint:  $F_{u,v} = \frac{R_1/u}{R_2/v}$ , where  $R_1 \sim \chi_u^2$ ,  $R_2 \sim \chi_v^2$ , and  $R_1$  and  $R_2$  are independent.]
  - (b) Let  $\mathbf{M} \sim W_p(k, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma}$  is positive definite and  $k > p + 1$ . Show that  $E[\mathbf{M}] = k\boldsymbol{\Sigma}$ .
  - (c) Let  $\mathbf{M}$  be as in part (b). Show that  $E[\mathbf{M}^{-1}] = c\boldsymbol{\Sigma}^{-1}$  for some constant  $c > 0$ , and find an explicit expression for this constant. [Hint: Use Property (5) of Wishart distribution given in [Section 12.2](#)].
- 12.2.** Consider a repeated measures study in which, for each of the  $n$  randomly selected subjects, an attribute (such as growth) is recorded at times  $t_1, \dots, t_p$ . Thus for the  $i$ th individual, the vector of measurements is  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{ip})^T$ . For each of the three popular models given below,  $\beta_0$  is a constant,  $\{\alpha_i\}$  are iid  $N(0, \sigma_\alpha^2)$ ,  $\{\delta_{ij}\}$  are iid  $N(0, \sigma_\delta^2)$ :
- (i)  $Y_{ij} = \beta_0 + \alpha_i + \tau_j + \delta_{ij}$ ,  $\tau_j$ 's are time effects (fixed) with  $\sum \tau_j = 0$ .
  - (ii)  $Y_{ij} = \beta_0 + \alpha_i + \beta_1 t_j + \delta_{ij}$ ,  $\beta_1$  is the slope (constant).
  - (iii)  $Y_{ij} = \beta_0 + \alpha_i + \beta_1 t_j + \delta_{ij}$ ,  $\{\beta_{i1}\}$  are iid random slopes with  $\beta_{i1} \sim N(\beta_1, \sigma_\beta^2)$ .
- It is understood that  $\{\alpha_i\}$ ,  $\{\beta_{i1}\}$ ,  $\{\delta_{ij}\}$  are all mutually independent. For each of the three models above,  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  are iid  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Explicitly obtain the elements of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  in each of the three cases.
- 12.3.** Let  $\mathbf{Y} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{u}$  is a vector in  $\mathbb{R}^p$ . Let  $Z = \mathbf{u}^T (\mathbf{Y} - \boldsymbol{\mu}) / \sqrt{\mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u}}$ .
- (a) If  $\mathbf{u}$  is nonrandom, then show that  $Z \sim N(0, 1)$ .
  - (b) Now assume that  $\mathbf{u}$  is a random vector which is independent of  $\mathbf{Y}$  and  $P[\mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} = 0] = 0$ . Show that  $Z \sim N(0, 1)$  and  $Z$  is independent of  $\mathbf{u}$ .
- 12.4.** Suppose that the growth of  $n$  randomly selected children are observed at  $p$  different time points and let  $\mathbf{Y}_i$  be the  $p$ -dim vector of observed growths for the  $i$ th child. Assume that  $\boldsymbol{\epsilon}_i = \mathbf{Y}_i - E[\mathbf{Y}_i]$  are iid  $N_p(\mathbf{0}, \boldsymbol{\Sigma})$ ,  $\boldsymbol{\Sigma}$  unknown. It is also

assumed that the mean growth at time  $t$  can be modeled by  $\beta_0 + \beta_1 t + \cdots + \beta_d t^d$ , where  $\{\beta_j\}$  are unknown and need to be estimated.

- (a) If growth of each of the  $n$  children are observed at times  $t_1, \dots, t_p$ , then  $E[Y_i] = \mathbf{D}\beta$  and obtain an explicit expression for the matrix  $\mathbf{D}$ . Explicitly write down  $-2 \times \log\text{-likelihood}$ .
- (b) If the observation time points for these children are not necessarily the same, that is, the growth of the  $i$ th child is measured at times  $t_{ij}, j = 1, \dots, p$ , then  $E[Y_i] = \mathbf{D}_i\beta$  and obtain an explicit expression for  $\mathbf{D}_i$ . Explicitly write down  $-2 \times \log\text{-likelihood}$ .
- (c) Show that the likelihood equations (ie, the minimizers of  $-2 \times \log\text{-likelihood}$ ) for part (a) are

$$\beta = (\mathbf{D}^T \Sigma^{-1} \mathbf{D})^{-1} \mathbf{D}^T \Sigma^{-1} \bar{\mathbf{Y}}, \quad \Sigma = n^{-1} \sum (Y_i - \mathbf{D}\beta)(Y_i - \mathbf{D}\beta)^T.$$

[Estimates of  $\beta$  and  $\Sigma$  are obtained by iterations starting with some reasonable estimates of  $\beta$  and  $\Sigma$ . For instance one may start with an initial estimate of  $\beta$  as  $\hat{\beta}^0 = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \bar{\mathbf{Y}}$  and use this in the expression of  $\Sigma$  (in the likelihood equations) to get an intial estimate  $\hat{\Sigma}^0$  of  $\Sigma$ .]

- (d) Show that the likelihood equations (ie, the minimizers of  $-2 \times \log\text{-likelihood}$ ) for part (b) are

$$\begin{aligned} \beta &= \left( \sum \mathbf{D}_i^T \Sigma^{-1} \mathbf{D}_i \right)^{-1} \left( \sum \mathbf{D}_i^T \Sigma^{-1} \bar{\mathbf{Y}} \right), \\ \Sigma &= n^{-1} \sum (Y_i - \mathbf{D}_i\beta)(Y_i - \mathbf{D}_i\beta)^T. \end{aligned}$$

- 12.5.** Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $N_p(\mu, \Sigma)$ , but assume that  $\Sigma$  is known. The statistic  $T_0^2 = n(\bar{\mathbf{Y}} - \mu_0)^T \Sigma^{-1} (\bar{\mathbf{Y}} - \mu_0)$  may be used for testing  $H_0: \mu = \mu_0$  against  $H_1: \mu \neq \mu_0$ .
- (a) Find the distribution of  $T_0^2$  under  $H_0$ . If  $H_0$  is not true, what is the distribution of  $T_0^2$ ?
  - (b) Show that  $E[T_0^2] = p + n(\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0)$ , when the true mean is  $\mu$ .
- 12.6.** Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $N_p(\mu, \Sigma)$ . Consider the problem of testing  $H_0: \mu = \mu_0$  against  $H_1: \mu \neq \mu_0$ . Then the Hotelling's  $T^2$  statistic is  $T^2 = n(\bar{\mathbf{Y}} - \mu_0)^T \mathbf{S}^{-1} (\bar{\mathbf{Y}} - \mu_0)$ , where  $\mathbf{S}$  is the sample covariance matrix. For  $\mathbf{a}$  in  $\mathbb{R}^p$ , we denote  $\mathbf{Y}(\mathbf{a}) = \mathbf{a}^T \mathbf{Y}$ ,  $\mu(\mathbf{a}) = \mathbf{a}^T \mu$ ,  $\varepsilon(\mathbf{a}) = \mathbf{a}^T \varepsilon$ ,  $\mu_0(\mathbf{a}) = \mathbf{a}^T \mu_0$ , and  $s(\mathbf{a})^2 = \mathbf{a}^T \mathbf{S} \mathbf{a}$ .
- (a) Consider the univariate model  $Y_i(\mathbf{a}) = \mu(\mathbf{a}) + \varepsilon_i(\mathbf{a})$ . Then for the problem of testing  $H_0: \mu(\mathbf{a}) = \mu_0(\mathbf{a})$  against  $H_1: \mu(\mathbf{a}) \neq \mu_0(\mathbf{a})$ , we may use the  $t$ -statistic  $t(\mathbf{a}) = \sqrt{n}(\bar{Y}(\mathbf{a}) - \mu_0(\mathbf{a}))/s(\bar{Y}(\mathbf{a}))$ , where  $s^2(\bar{Y}(\mathbf{a})) = \mathbf{a}^T \mathbf{S} \mathbf{a}/n$ . Show that  $\max_{\mathbf{a} \neq 0} t(\mathbf{a})^2 = T^2$ .
  - (b) Show that  $E[T^2] = \frac{p(n-1)}{n-p-2}$  when  $H_0$  is true.
  - (c) Show that  $E[T^2] = \frac{n(n-1)}{n-p-2}(\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0) + \frac{p(n-1)}{n-p-2}$ . Show that  $E[T^2] = \frac{p(n-1)}{n-p-2}$  when and only when  $H_0$  is true.

- 12.7.** Let  $Y_1, \dots, Y_n$  be iid  $N_p(\mu, \Sigma)$ . Let  $\theta_i = \mu_{i+2} - 2\mu_{i+1} + \mu_i$ ,  $i = 1, \dots, p-2$ . Suppose that we are interested in testing  $H_0: \theta = \mathbf{0}$  against  $H_1: \theta \neq \mathbf{0}$ , where  $\theta$  is the  $(p-2)$ -dim vector consisting of elements  $\theta_1, \dots, \theta_{p-2}$ . [This kind of testing may be important in growth analysis where we may want to test if  $\mu_i$ 's are linear in  $i$ .]
- (a) Express  $\theta$  as a linear function  $D\mu$  of  $\mu$  by finding a matrix  $D$  explicitly.
  - (b) Find the MLE  $\hat{\theta}$  of  $\theta$ . Obtain the distribution of this estimate. Obtain an unbiased estimate of the covariance matrix of  $\hat{\theta}$ .
  - (c) Obtain an appropriate statistic for testing  $H_0$  against  $H_1$ , and find the distribution of this statistic under  $H_0$ .
- 12.8.** Let  $Y_1, \dots, Y_n$  be iid  $N_{2p}(\mu, \Sigma)$ . Let  $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$  where  $\mu_1$  and  $\mu_2$  are  $p$ -dim. We are interested in inference on  $\theta = \mu_1 - \mu_2$ . This type of issue comes up in the case of paired observations. Assume that  $\Sigma$  is of the form  $\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$ , where  $\Sigma_{11}, \Sigma_{12}, \Sigma_{21}$ , and  $\Sigma_{22}$  are  $p \times p$  matrices. Similarly the sample covariance matrix  $S$  can be written as  $\begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}$ , where  $S_{11}, S_{12}, S_{21}$ , and  $S_{22}$  are  $p \times p$  matrices.
- (a) Express  $\theta$  in the form  $D\mu$  and find the matrix  $D$  explicitly. Find the MLE of  $\theta$  and then find the distribution of this estimate. Obtain the parameters of this distribution explicitly in terms of  $\theta$ ,  $\Sigma_{11}$ ,  $\Sigma_{12}$ ,  $\Sigma_{21}$ , and  $\Sigma_{22}$ .
  - (b) Suppose you want to test  $H_0: \theta = \mathbf{0}$  against  $H_1: \theta \neq \mathbf{0}$ . Find the appropriate Hotelling's  $T^2$ -statistic and obtain its distribution under  $H_0$ . Modify this test statistic appropriately if  $\Sigma$  were known. Find the distribution of this modified statistic under  $H_0$ .
  - (c) Obtain simultaneous confidence intervals for  $\theta_1, \dots, \theta_m$  with family confidence of at least  $1 - \alpha$  using the Bonferroni approach. Obtain an explicit expression for the standard error of the estimate of each  $\theta_i$  as a combination of the elements of the matrices  $S_{11}, S_{12}, S_{21}$ , and  $S_{22}$ .
  - (d) In some cases it is of interest to test  $H_0: \theta_1 + \dots + \theta_m \leq 0$  against  $H_1: \theta_1 + \dots + \theta_m > 0$  at a level of significance  $\alpha$ , where  $m \leq p$ . Obtain an appropriate statistic for this testing problem and state the decision rule.
- 12.9.** Consider a one-factor MANOVA model  $Y_{ij} = \mu_i + \epsilon_{ij}$ ,  $j = 1, \dots, n_i$ ,  $i = 1, \dots, k$ , where  $\{\epsilon_{ij}\}$  are iid  $N_p(\mathbf{0}, \Sigma)$ . Let  $B$ ,  $W$ , and  $T$  denote the between group, within group, and total SSP matrices, respectively.
- (a) Show that the matrices  $B$ ,  $W$ , and  $T$  are nonnegative definite.
  - (b) Show that  $|W| \leq |T|$ .
  - (c) Assume that  $T$  is positive definite with probability 1. Show that  $|W| = |T|$  when and only when  $\bar{Y}_i = \bar{Y}_{..}$  for  $i = 1, \dots, k$ . [Hint for parts (b) and (c): Look at  $T^{-1/2}WT^{-1/2}$ .]
  - (d) Let  $\mu = \sum(n_i/n)\mu_i$  and  $\alpha_i = \mu_i - \mu$ , where  $n = n_1 + \dots + n_k$ . Let  $\hat{\alpha}_i = \bar{Y}_i - \bar{Y}_{..}$  Show that  $\bar{Y}_{..} \sim N_p(\mu, n^{-1}\Sigma)$  and  $\hat{\alpha}_i \sim N_p(\alpha_i, (n_i^{-1} - n^{-1})\Sigma)$ .

- (e) Let  $\theta = \sum c_i \alpha_i$  be a contrast, that is,  $\{c_i\}$  are real numbers and they satisfy the constraint  $\sum c_i = 0$ . Let  $\hat{\theta} = \sum c_i \hat{\alpha}_i$ . Show that  $\hat{\theta} \sim N_p(\theta, \sum c_i^2/n_i \Sigma)$ .
- (f) First prove that  $E[\mathbf{B}] = (k-1)\Sigma + \sum m_i \alpha_i \alpha_i^T$ . Then show that  $E[\mathbf{B}] = (k-1)\Sigma$  when and only when  $\mu_1 = \dots = \mu_k$ .
- 12.10.** Consider a one-factor MANOVA model as in [Exercise 12.9](#). Consider the transformed data  $Z_{ij} = \mathbf{a} + \mathbf{A}Y_{ij}$ , where  $\mathbf{a}$  is a vector in  $\mathbb{R}^p$  and  $\mathbf{A}$  is a nonsingular matrix of order  $p \times p$ .
- Write down a one-factor MANOVA model for  $\{Z_{ij}\}$ , that is,  $Z_{ij} = \theta_i + \delta_{ij}$ , where  $\delta_{ij}$ 's have zero means. Find the distribution of  $\{\delta_{ij}\}$ . Express  $\theta_i$  and the parameters of the distribution of  $\delta_{ij}$  explicitly in terms of  $\mu_i$ ,  $\mathbf{a}$ ,  $\mathbf{A}$ , and  $\Sigma$ .
  - Show that the between group, within group, and total SSP matrices for the transformed data  $\{Z_{ij}\}$  are  $\mathbf{ABA}^T$ ,  $\mathbf{AWA}^T$ , and  $\mathbf{ATA}^T$ , respectively.
  - Consider the problem of testing  $H_0: \theta_1 = \dots = \theta_k$  against  $H_1$ : not all  $\theta_i$ 's are the same. Show that the Wilks' lambda, Pillai trace, and Roy's largest root statistics obtained on the basis of the data  $\{Z_{ij}\}$  are the same as those for  $\{Y_{ij}\}$ .
  - When  $k = 2$ , show that the Lawley-Hotelling statistic  $\text{trace}(\mathbf{BW}^{-1})$  is equal to  $cT^2$ , for some constant  $c > 0$ , where  $T^2$  is the two-sample  $T^2$ -statistic for testing  $H_0: \mu_1 = \mu_2$  vs  $H_1: \mu_1 \neq \mu_2$ . Find the constant  $c$ .
- 12.11.** Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be iid  $p$ -dim random vectors with mean vector  $\mu$  and covariance matrix  $\Sigma$  whose eigenvalues are  $\lambda_1 \geq \lambda_2 \geq \dots$  with the corresponding orthonormal eigenvectors  $\mathbf{u}_1, \mathbf{u}_2, \dots$ . Let  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$  be the eigenvalues of the sample covariance matrix  $\mathbf{S}$  with the corresponding orthonormal eigenvectors  $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$ . The sample values of the  $j$ th principal component are  $\hat{Z}_{tj} = \hat{\mathbf{u}}_j^T (\mathbf{Y}_t - \bar{\mathbf{Y}})$ ,  $t = 1, \dots, n$ .
- Show that the sample mean and sample variance of  $\{\hat{Z}_{tj}: t = 1, \dots, n\}$  are zero and  $\hat{\lambda}_j$ , respectively.
  - Show that the sample correlation of  $\{\hat{Z}_{ti}, \hat{Z}_{tj}\}: t = 1, \dots, n$ ,  $i \neq j$ , is equal to zero.
- For parts (c) and (d) assume that the population is normal and that the eigenvalues  $\lambda_1, \dots, \lambda_p$  are distinct.
- Let  $\pi_k = (\lambda_1 + \dots + \lambda_k)/\text{trace}(\Sigma)$  and  $\hat{\pi}_k = (\hat{\lambda}_1 + \dots + \hat{\lambda}_k)/\text{trace}(\mathbf{S})$ ,  $k < p$ . Let  $\Sigma_k = \sum_{j=1}^k \lambda_j \mathbf{u}_j \mathbf{u}_j^T$ . Show that  $\sqrt{n}(\hat{\pi}_k - \pi_k) \xrightarrow{D} N(0, 2\tau_k^2)$ , where  $\tau_k^2 = [(1 - \pi_k)^2 \text{trace}(\Sigma_k^2) + \pi_k^2 \text{trace}((\Sigma - \Sigma_k)^2)]/\text{trace}(\Sigma)^2$ .
  - Let  $\theta$  and  $\hat{\theta}$  be the geometric means of  $\{\lambda_1, \dots, \lambda_p\}$  and  $\{\hat{\lambda}_1, \dots, \hat{\lambda}_p\}$ , respectively. Show that  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{D} N(0, 2\theta^2/p)$ .
- 12.12.** Consider the factor model given in Eq. (3) and assume that  $\mathbf{L}$  and  $\Psi$  are known and  $\mu = \mathbf{0}$ . Suppose it is desired to predict the vector of common factors  $\mathbf{f}$  when  $\mathbf{Y}$  is given and this problem examines the three predictors  $\hat{\mathbf{f}}^{(1)}, \hat{\mathbf{f}}^{(2)}$ , and  $\hat{\mathbf{f}}^{(3)}$  given in [Section 12.10.2](#).
- Show that  $\hat{\mathbf{f}}^{(3)}$  is the best linear unbiased predictor of  $\mathbf{f}$ .

- (b) For any  $\mathbf{a} \in \mathbb{R}^k$ , explicitly calculate  $E\left[\mathbf{a}^T \hat{\mathbf{f}}^{(1)} - \mathbf{a}^T \mathbf{f}\right]^2$ ,  $E\left[\mathbf{a}^T \hat{\mathbf{f}}^{(2)} - \mathbf{a}^T \mathbf{f}\right]^2$ , and  $E\left[\mathbf{a}^T \hat{\mathbf{f}}^{(3)} - \mathbf{a}^T \mathbf{f}\right]^2$ , and rank these three predictors  $\hat{\mathbf{f}}^{(1)}$ ,  $\hat{\mathbf{f}}^{(2)}$ , and  $\hat{\mathbf{f}}^{(3)}$  from the worst to the best.
- 12.13.** Let  $\mathbf{Y} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and assume that the factor model given in Eq. (3) holds, that is,  $\mathbf{Y} = \boldsymbol{\mu} + \mathbf{L}\mathbf{f} + \boldsymbol{\varepsilon}$ , where  $\mathbf{L}$  is  $p \times k$ ,  $k < p$ ,  $\mathbf{f} \sim N_k(\mathbf{0}, I)$ ,  $\boldsymbol{\varepsilon} \sim N_p(\mathbf{0}, \boldsymbol{\Psi})$ , where  $\boldsymbol{\Psi}$  is diagonal, and  $\mathbf{f}$  and  $\boldsymbol{\varepsilon}$  are assumed to be independent. Suppose it is desired to predict  $\mathbf{f}$  when  $\mathbf{Y}$  is given and the matrices  $\mathbf{L}$  and  $\boldsymbol{\Psi}$  are known and  $\boldsymbol{\mu} = \mathbf{0}$ .
- (a) Find the conditional distribution of  $\mathbf{f}$  given  $\mathbf{Y}$  and show that the mean of this conditional distribution is  $\hat{\mathbf{f}}^{(3)}$  as given in Section 12.10.2.
  - (b) Using the joint distribution of  $\mathbf{Y}$  and  $\mathbf{f}$ , obtain an analog of the mixed model equations given in Section 11.10.2 of Chapter 11 (assuming that  $\mathbf{L}$  and  $\boldsymbol{\Psi}$  are known and  $\boldsymbol{\mu} = \mathbf{0}$ ) and show that  $\hat{\mathbf{f}}^{(3)}$  is the solution of this equation.
- 12.14.** Let  $\mathbf{Y}$  be from one of the populations  $N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$  (population 1) or  $N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$  (population 2) with probabilities  $\pi_1$  and  $\pi_2$ , respectively, with  $\pi_1 + \pi_2 = 1$ . Consider a linear discriminant rule with prior probabilities  $\pi_1$  and  $\pi_2$ . Assume that  $\boldsymbol{\mu}_1$ ,  $\boldsymbol{\mu}_2$ ,  $\boldsymbol{\Sigma}$ , and  $\pi_1$  are known. Denote  $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  and  $\log(\pi_1/\pi_2)$  by  $\delta^2$  and  $c$ , respectively.
- (a) Let  $R = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \bar{\boldsymbol{\mu}})$ , where  $\bar{\boldsymbol{\mu}} = (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)/2$ . Show that the linear discriminant rule allocates  $\mathbf{Y}$  to population 1 if  $R > -c$ , and allocates  $\mathbf{Y}$  to population 2 if the inequality is reversed.
  - (b) Show that  $R$  is  $N(\delta^2/2, \delta^2)$  if  $\mathbf{Y}$  is from population 1, and  $R$  is  $N(-\delta^2/2, \delta^2)$  if  $\mathbf{Y}$  is from population 2.
  - (c) Show that the probability of misclassification is equal to  $\Phi(-\delta/2 - c/\delta)$  if  $\mathbf{Y}$  is from population 1, and is equal to  $\Phi(-\delta/2 + c/\delta)$  when  $\mathbf{Y}$  is from population 2. [Here  $\Phi$  is the cdf of the standard normal distribution.]
  - (d) If  $\mathbf{Y}$  is from one of the two populations (with probability  $\pi_1$  from population 1 and with probability  $\pi_2$  from population 2) and it is classified using a linear discriminant rule, find the probability of misclassifying  $\mathbf{Y}$ , in terms of  $\delta$ ,  $\pi_1$ , and  $\pi_2$ .
- 12.15.** Prove Theorem 12.12.2.
- 12.16.** Prove Lemma 12.12.1.

# Time Series

## 13.1 Introduction

In the previous chapters, except for random-effects models in [Chapter 11](#), all the different types of statistical modeling and procedures are concerned with data sets consisting of independent observations. In practice, however, there are many cases when the assumption of independence is not tenable and this is particularly true when the observations are recorded over time and/or over space (geographical locations). A simple example for such a data set is daily records of average levels of ozone concentration in the air at various locations in a particular geographical region. This chapter is only concerned with data sets consisting of observations recorded over time and such observations are usually dependent. For instance, when unemployment rates are recorded over months, observation in a particular month depends on the employment levels in the previous months. Methodologies developed for the investigation of such data sets are called time series methods. Here are a few examples of time series data:

- (a) Annual precipitation at Lake Michigan in the last 75 years.
- (b) Annual temperature anomalies (ie, average yearly temperature minus a base value) in the last 150 years.
- (c) Monthly unemployment in the United States in the last 50 years.
- (d) Monthly electricity sales to the residential sector in the United States in the last 50 years.
- (e) EEG data (used in diagnosing patients) at a particular skull location.

We may denote the observations as  $\{Y_t: t = 1, \dots, n\}$  where the time unit may be a year or a month or a week or even fraction of a second depending on the particular problem at hand. An actual examination of the data in the first example indicates that the annual precipitation  $\{Y_t: t = 1, \dots, n = 75\}$  tends to fluctuate around a constant value. In the second example, there is an overall increase in the annual temperature anomalies  $\{Y_t: t = 1, \dots, n = 150\}$  over the last 150 years and thus it is reasonable to model  $Y_t$  as a smooth part (trend) plus a rough part (random errors which are identically distributed, but not necessarily independent). An examination of the monthly electricity sales series  $\{Y_t: t = 1, \dots, n = 12 \times 50 = 600\}$  would reveal an overall increase over time. Moreover,

there is a seasonal factor (ie, January sales in consecutive years tend to be similar, February sales in consecutive years tend to be similar, and so on). There may also be a cyclical factor (about 9–12 years) as the sales may dip a little after economic recessions. Thus  $Y_t$  may be modeled as a sum of the trend, the seasonal effect, the cyclical effect, and a mean zero random error. Thus a general model for  $Y_t$  for three of the five examples given above may be

$$\text{Example (a): } Y_t = \mu + X_t,$$

$$\text{Example (b): } Y_t = \mu_t + X_t,$$

$$\text{Example (d): } \log Y_t = \mu_t + S_t + C_t + X_t,$$

where  $\mu_t$  is the trend,  $\{S_t\}$  are the seasonal effects,  $\{C_t\}$  are the cycles, and  $\{X_t\}$  are the random errors which are identically distributed but may not be independent. Note that in Example (d), the series is transformed by the natural logarithm in order achieve constant variability over time and this will be discussed below. The trend can be modeled as a smooth function by nonparametric methods. Seasonals  $\{S_t\}$  may be modeled by various methods including linear combination of sines and cosines. If the period of cycles are known, they may also be modeled by sines and cosines. Apart from the above-mentioned methods for modeling the trend, the seasonals, and the cycles, there are also probabilistic modeling schemes. Forecasting is one of the important goals in time series analysis. Forecasting at time  $t + h$ ,  $h \geq 1$ , requires estimates of  $\mu_{t+h}$ ,  $S_{t+h}$ ,  $C_{t+h}$ , and  $X_{t+h}$  which can be added to get a forecast value of  $Y_{t+h}$  as  $\hat{Y}_{t+h} = \hat{\mu}_{t+h} + \hat{S}_{t+h} + \hat{C}_{t+h} + \hat{X}_{t+h}$ .

In the fifth example above, the observations are combination of waves (alpha, beta, theta, delta, etc.), and the weights of the combinations vary depending on whether the subject is normal or has a disease such as epilepsy. For instance, theta waves are in the frequency range of 3–8 Hz and are present in diseased patients, whereas alpha waves vary in the frequency range of 8–13 Hz and are present in normal individuals without any external stimulus. The goal of the analysis here is not prediction, but to determine the combination of different waves in the observed series for a particular patient and this is done by using the spectral analysis which is described in the last section of this chapter.

This chapter is concerned with understanding and analysis of the rough part  $\{X_t\}$ . The main assumptions are

- (i)  $\{X_t\}$  are identically distributed,
- (ii) for any  $t$ , the correlation between  $X_t$  and  $X_{t+h}$ ,  $h \geq 0$ , depends only on  $h$ , and
- (iii) the correlation between  $X_t$  and  $X_{t+h}$  is negligible when  $h$  is large.

A series with properties (i) and (ii) is called *stationary*, and property (iii) is a statistical necessity, which allows forecasting and consistent estimation for parametric models based on observed data sets.

### *Types of Nonstationary Series*

There may be many sources of nonstationarity, but we briefly point out three of them:

- (i) unequal means over time,
- (ii) unequal variance over time, and
- (iii) inappropriate time scale.

Often, a reasonable way to view unequal means is to treat it as a smooth function of time (trend). For some series, such as daily records of sulfur dioxide levels in the air, sharp changes may correspond to unusual events such as volcano eruptions, and may need to be incorporated into the trend.

In the analysis of financial time series data such as rate of return on stocks, the main focus of investigation is on modeling the variance which changes over time. In other cases, unequal variance can sometimes be remedied by transforming the data using a Box-Cox transformation

$$Y_t(\lambda) = \begin{cases} \frac{Y_t^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \log(Y_t) & \lambda = 0. \end{cases}$$

If a transformation turns out to be useful, then analysis and modeling are done on the transformed series. For instance, a logarithmic transformation for monthly electricity sales (Example (d) above) is appropriate for achieving equal variance.

For some series, such as a signal from bird chirping, attempts are made to plot, understand, and analyze the data on a time scale  $H(t)$  (a nonlinear function of time  $t$ ) on which the series may be stationary. Typically  $H$  is unknown and needs to be estimated.

### *A Simple Method for Extracting the Stationary Part*

Consider a series of the form  $Y_t = m_t + S_t + X_t$ , where  $\{m_t\}$  is the smooth trend,  $\{S_t\}$  are the seasonal effects with period  $s$ , and  $\{X_t\}$  is the stationary part. As mentioned above, in some cases approximate equal variance may be achieved by employing Box-Cox transformations. When the seasonal effects are absent (ie, the model is  $Y_t = m_t + X_t$ ), a popular method for the analysis of such series employs the integrated autoregressive-moving average (ARIMA) models. When both the trend and the seasonals are present, then such a series can be analyzed by using what is known as the integrated seasonal autoregressive-moving average (seasonal ARIMA) model. A good detailed description of these models can be found in the book by Box et al. [62].

Here we briefly describe a simple regression method for estimating  $m_t$  and  $S_t$ , and this procedure works well for some series. Since  $X_t = Y_t - (m_t + S_t)$ , we can obtain an approximation of  $X_t$  if we can obtain estimates of  $m_t$  and  $S_t$ . For the sake of identifiability, let us assume that  $S_{t-r+1} + \dots + S_t = 0$  for any  $t$ , where  $r$  is the seasonal order. In order to simplify the discussion, let us assume that we are dealing with quarterly data (ie,  $r = 4$ ). If we model the trend  $m_t$  by a polynomial of degree  $d$ , then

$$m_t = \beta_0 + \beta_1 t + \cdots + \beta_d t^d.$$

In order to account for the seasonal effects, create variables,  $I_{t1}$ ,  $I_{t2}$ , and  $I_{t3}$  as

$$I_{t1} = \begin{cases} 1 & \text{if time } t \text{ is Quarter 1} \\ -1 & \text{if time } t \text{ is Quarter 4} \\ 0 & \text{otherwise} \end{cases}, \quad I_{t2} = \begin{cases} 1 & \text{if time } t \text{ is Quarter 2} \\ -1 & \text{if time } t \text{ is Quarter 4} \\ 0 & \text{otherwise} \end{cases},$$

$$I_{t3} = \begin{cases} 1 & \text{if time } t \text{ is Quarter 3} \\ -1 & \text{if time } t \text{ is Quarter 4} \\ 0 & \text{otherwise.} \end{cases}$$

Thus  $m_t + S_t$  can be modeled as

$$\beta_0 + \beta_1 t + \cdots + \beta_d t^d + \theta_1 I_{t1} + \theta_2 I_{t2} + \theta_3 I_{t3},$$

and the unknown parameters  $\beta_0, \dots, \beta_d$  and  $\theta_1, \theta_2$ , and  $\theta_3$  can be estimated by minimizing the least squares criterion

$$\sum_{t=1}^n [Y_t - \beta_0 - \beta_1 t - \cdots - \beta_d t^d - \theta_1 I_{t1} - \theta_2 I_{t2} - \theta_3 I_{t3}]^2.$$

Once the estimates of these parameters are obtained, one can get estimates  $\hat{m}_t$  and  $\hat{S}_t$  of  $m_t$  and  $S_t$ . Then an estimate of  $X_t$  is given by  $\hat{X}_t = Y_t - (\hat{m}_t + \hat{S}_t)$ .

It is important to note that the seasonal fluctuations  $\{S_t\}$  are assumed to be nonrandom in this discussion, which may not always be appropriate.

## 13.2 Concept of Stationarity

In order to develop an appropriate mathematical framework, it is usually assumed that the data  $\{X_1, \dots, X_n\}$  is a finite section of an infinite series  $\{X_t: -\infty < t < \infty\}$ . Throughout this chapter, we use the terms “series” and “process” to mean a sequence of random variables. A series  $\{X_t\}$  is called *strictly stationary* if  $\{X_t, \dots, X_{t+k}\}$  has the same distribution as  $\{X_{t+h}, \dots, X_{t+h+k}\}$  for any  $t$  and  $h \geq 0$ . This notion of stationarity is usually too strict to be useful in applications since it is difficult to verify it in practice. A weaker version, also known as covariance stationarity, assumes that  $\text{Cov}[X_t, X_{t+h}]$  depend only on  $h$ . It is easy to see that strict stationarity implies weak stationarity. From now on, we assume that  $\{X_t\}$  has mean  $\mu$  and  $\gamma(h) = \text{Cov}[X_t, X_{t+h}]$  for any  $t$  and  $h$ . Since  $\gamma(-h) = \text{Cov}[X_{t+h}, X_t]$  and  $\text{Cov}[X_t, X_{t+h}] = \text{Cov}[X_{t+h}, X_t]$ , it thus follows that  $\gamma(-h) = \gamma(h)$  for any integer  $h$ . The correlation between  $X_t$  and  $X_{t+h}$  can be easily seen to be  $\rho(h) = \gamma(h)/\gamma(0)$  and this quantity also does not depend on  $t$ . The sequences  $\{\gamma(h)\}$  and  $\{\rho(h)\}$  are called the autocovariance function and the autocorrelation function, respectively.

A time series is *stationary Gaussian* if the joint distribution of  $\{X_{t+1}, \dots, X_{t+p}\}$ , for any  $t$  and  $p \geq 1$ , has a  $p$ -dim multivariate normal distribution with the mean vector  $(\mu, \dots, \mu)^T$

and covariance matrix  $\boldsymbol{\Gamma}_p$  with elements  $\gamma(k - j)$ ,  $1 \leq j, k \leq p$ . It can be shown that a covariance-stationary Gaussian time series is also “strictly” stationary.

**Definition 13.2.1.** A series  $\{X_t\}$  is called (weakly) stationary if, for any  $t$ ,  $E[X_t] = \mu$  and  $\text{Cov}[X_t, X_{t+h}]$  depends only on  $h$ . Its autocovariance  $\{\gamma(h)\}$  and autocorrelation functions  $\{\rho(h)\}$  are  $\gamma(h) = \text{Cov}[X_t, X_{t+h}]$  and  $\rho(h) = \text{Corr}[X_t, X_{t+h}] = \gamma(h)/\gamma(0)$ .

We now present a few examples of stationary series and one example of a nonstationary series.

**Example 13.2.1** (White Noise). Let  $\{X_t\}$  be iid with mean  $\mu$  and variance  $\sigma^2$ , then it is stationary with

$$\gamma(h) = \begin{cases} \sigma^2 & \text{if } h = 0 \\ 0 & \text{if } h \neq 0 \end{cases} \text{ and } \rho(h) = \begin{cases} 1 & \text{if } h = 0 \\ 0 & \text{if } h \neq 0 \end{cases}.$$

If  $\{X_t\}$  are iid with  $\mu = 0$ , then such a series is called white noise. Some authors dispense with the iid assumption and call a series white noise if  $\{X_t\}$  have zero mean, identical variances, and are mutually uncorrelated.

**Example 13.2.2** (Moving Average Process). If a series  $\{X_t\}$  can be written as

$$X_t - \mu = \varepsilon_t + \theta \varepsilon_{t-1},$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ , then it is called a moving average process of order 1. For such a series,  $E[X_t] = \mu$ , and the autocovariances are

$$\begin{aligned} \gamma(0) &= \text{Var}[X_t] = \text{Var}[\varepsilon_t] + \text{Var}[\theta \varepsilon_{t-1}] = (1 + \theta^2)\sigma^2, \\ \gamma(1) &= \text{Cov}[X_t, X_{t-1}] = \text{Cov}[\varepsilon_t + \theta \varepsilon_{t-1}, \varepsilon_{t-1} + \theta \varepsilon_{t-2}] \\ &= \text{Cov}[\theta \varepsilon_{t-1}, \varepsilon_{t-1}] = \theta \sigma^2, \text{ and} \\ \gamma(h) &= 0, \quad h \geq 2. \end{aligned}$$

Clearly, the autocorrelations are

$$\rho(1) = \theta/(1 + \theta^2), \quad \rho(h) = 0, \quad h \geq 2.$$

A process more general than the above is

$$X_t - \mu = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \cdots + \theta_q \varepsilon_{t-q},$$

where  $\{\varepsilon_t\}$  are mean 0 iid variables with common variance  $\sigma^2$ . This is known as a moving average process of order  $q$  and is denoted by  $MA(q)$ . For this series,  $E[X_t] = \mu$ , and the autocovariances and autocorrelations are

$$\begin{aligned} \gamma(h) &= \sigma^2 \sum_{j=0}^{q-h} \theta_j \theta_{j+h}, \quad 0 \leq h \leq q, \text{ with } \theta_0 = 1, \\ \gamma(h) &= 0, \quad h \geq q+1, \\ \rho(h) &= \sum_{j=0}^{q-h} \theta_j \theta_{j+h} / \left( \sum_{j=0}^q \theta_j^2 \right), \quad 0 \leq h \leq q, \text{ and} \\ \rho(h) &= 0, \quad h \geq q+1. \end{aligned}$$

**Example 13.2.3** (Autoregressive Process). If a series  $\{X_t\}$  is representable as

$$X_t - \mu = \phi(X_{t-1} - \mu) + \varepsilon_t,$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ , then it is called an autoregressive process of order 1. For this series,  $E[X_t] = \mu$ . In order to obtain the variance of  $X_t$ , note that  $Cov[X_{t-1}, \varepsilon_t] = 0$ , and hence

$$\begin{aligned}\gamma(0) &= \text{Var}[X_t] = \text{Var}[\phi X_{t-1} + \varepsilon_t] = \phi^2 \gamma(0) + \sigma^2, \text{ and} \\ \gamma(0) &= (1 - \phi^2)^{-1} \sigma^2.\end{aligned}$$

The condition  $\gamma(0) = \text{Var}[X_t] > 0$  requires that  $\phi^2 < 1$  (ie,  $|\phi| < 1$ ) and this will be shown later to be the condition for the series to be stationary. Autocovariances of this series can be obtained by using a recursive procedure. Since  $X_{t-1}$  is uncorrelated with  $\varepsilon_t$ , we have

$$\begin{aligned}\gamma(1) &= \text{Cov}[X_t, X_{t-1}] = \text{Cov}[\phi X_{t-1} + \varepsilon_t, X_{t-1}] \\ &= \phi \text{Cov}[X_{t-1}, X_{t-1}] = \phi \gamma(0).\end{aligned}$$

Similarly, noting that  $\varepsilon_t$  is uncorrelated with  $X_{t-h}$  for  $h \geq 1$ , we have

$$\begin{aligned}\gamma(h) &= \text{Cov}[X_t, X_{t-h}] = \text{Cov}[\phi X_{t-1} + \varepsilon_t, X_{t-h}] \\ &= \phi \text{Cov}[X_{t-1}, X_{t-h}] = \phi \gamma(h-1).\end{aligned}$$

Thus for  $h = 2, \dots$ , we have

$$\begin{aligned}\gamma(2) &= \phi \gamma(1) = \phi^2 \gamma(0), \\ \gamma(3) &= \phi \gamma(2) = \phi^3 \gamma(0), \dots\end{aligned}$$

This argument shows that for any  $h \geq 0$

$$\gamma(h) = \phi^h \gamma(0) \text{ with } \gamma(0) = (1 - \phi^2)^{-1} \sigma^2, \text{ and } \rho(h) = \phi^h.$$

Since  $|\phi| < 1$ ,  $\rho(h)$  converges to zero exponentially as  $h \rightarrow \infty$ .

A general version of the simple autoregressive process is

$$X_t - \mu = \phi_1(X_{t-1} - \mu) + \cdots + \phi_p(X_{t-p} - \mu) + \varepsilon_t,$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ . This is an autoregressive process of order  $p$  with mean  $E[X_t] = \mu$  and it is denoted by  $AR(p)$ . Conditions on  $\phi_1, \dots, \phi_p$  needed for the process  $\{X_t\}$  to be stationary will be discussed later. In the  $AR(p)$  case, it is not possible to obtain explicit expressions for the autocovariance or autocorrelation functions though computer packages such as R can be used to obtain them for given values of  $\phi_1, \dots, \phi_p$  and  $\sigma^2$ . It should be pointed out that, under the condition of stationarity, the autocorrelation function  $\rho(h)$  converges rapidly (exponentially) to zero as  $h \rightarrow \infty$ .

**Example 13.2.4** (Autoregressive-Moving Average Process). A series  $\{X_t\}$  which has both  $AR(p)$  and  $MA(q)$  parts is called an  $ARMA(p, q)$  process, and is representable in the form

$$\begin{aligned} X_t - \mu &= \phi_1(X_{t-1} - \mu) + \cdots + \phi_p(X_{t-p} - \mu) \\ &\quad + \varepsilon_t + \theta_1\varepsilon_{t-1} + \cdots + \theta_q\varepsilon_{t-q}, \end{aligned}$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ . Conditions on  $\phi_1, \dots, \phi_p$  are needed to guarantee stationarity as in the  $AR(p)$  case. In general there are no explicit forms for the autocovariances and autocorrelations even though computer packages can be used to obtain them. As in the  $AR(p)$  case, under appropriate conditions, the autocorrelation  $\rho(h)$  converges to zero exponentially as  $h \rightarrow \infty$ . The error terms  $\{\varepsilon_t\}$  in AR, MA, or ARMA series are sometimes called *innovations*.

**Example 13.2.5** (Random Walk). A mean zero  $AR(1)$  series with  $\phi = 1$  is called a random walk and it has the form  $X_t = X_{t-1} + \varepsilon_t$ ,  $t = 1, 2, \dots$ , where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ . It is easy to see that  $X_t = X_0 + \varepsilon_1 + \cdots + \varepsilon_t$  and assuming that  $E[X_0] = \mu$ , we have

$$E[X_t] = \mu \text{ and } \text{Cov}[X_t, X_{t+h}] = \text{Var}[X_0] + t\sigma^2, \quad h \geq 0.$$

Clearly this series is not stationary since  $\text{Cov}[X_t, X_{t+h}]$  depends on  $t$ . Random walks are sometimes used for modeling the trend of a nonstationary time series.

### 13.2.1 Representation of the Autocovariance Function

We begin with an important result on the representation of the covariance function  $\{\gamma(h) : -\infty < h < \infty\}$  of a stationary series. Under the assumption that  $\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty$ , the function

$$f(w) = \sum_{h=-\infty}^{\infty} \exp[-2\pi ihw]\gamma(h), \quad h \geq 0, \quad (1)$$

is well defined on  $[-1/2, 1/2]$ , where  $i = \sqrt{-1}$ . This function  $f$  is called the spectral density function of the series  $\{X_t\}$  with the autocovariance function  $\{\gamma(h)\}$ . It is periodic with period 1 (ie,  $f$  can be defined for any real  $w$  and  $f(w) = f(w+1)$ ). More detailed discussion on spectral density function will be given in a later section.

Since  $\gamma(h) = \gamma(-h)$  and  $\exp[-2\pi ihw] + \exp[2\pi ihw] = 2 \cos(2\pi hw)$ , we have

$$\begin{aligned} f(w) &= \sum_{h=-\infty}^{-1} \exp[-2\pi ihw]\gamma(h) + \gamma(0) + \sum_{h=1}^{\infty} \exp[-2\pi ihw]\gamma(h) \\ &= \gamma(0) + \sum_{h=1}^{\infty} \{\exp[2\pi ihw] + \exp[-2\pi ihw]\}\gamma(h) \\ &= \gamma(0) + 2 \sum_{h=1}^{\infty} \cos(2\pi hw)\gamma(h). \end{aligned}$$

Thus the spectral density function is a trigonometric series with autocovariances  $\{\gamma(h)\}$  as the coefficients. It is symmetric about 0 (ie,  $f(-w) = f(w)$ ), and is nonnegative (to be shown below). What is the connection between the autocovariance function

$\{\gamma(h)\}$  and the spectral density function? If we know the spectral density function  $f$ , then we can get the autocovariances as

$$\gamma(h) = \int_{-1/2}^{1/2} \exp[2\pi i h w] f(w) dw, \quad h \geq 0.$$

In order to verify that  $f$  is nonnegative, consider the arguments given below. The discrete cosine and sine transforms of  $\{X_t - \mu, t = 1, \dots, n\}$  are

$$\begin{aligned}\tilde{X}_{c,n}(w) &= n^{-1/2} \sum_{t=1}^n (X_t - \mu) \cos(2\pi w t), \text{ and} \\ \tilde{X}_{s,n}(w) &= n^{-1/2} \sum_{t=1}^n (X_t - \mu) \sin(2\pi w t),\end{aligned}$$

where  $|w| \leq 1/2$ . Then

$$\begin{aligned}\mathrm{E}[\tilde{X}_{c,n}(w)]^2 + \mathrm{E}[\tilde{X}_{s,n}(w)]^2 &= n^{-1} \sum_{1 \leq s, t \leq n} \cos(2\pi ws) \cos(2\pi wt) \mathrm{Cov}[X_s, X_t] \\ &\quad + n^{-1} \sum_{1 \leq s, t \leq n} \sin(2\pi ws) \sin(2\pi wt) \mathrm{Cov}[X_s, X_t] \\ &= n^{-1} \sum_{1 \leq s, t \leq n} \cos(2\pi w(t-s)w) \mathrm{Cov}[X_s, X_t] \\ &= n^{-1} \sum_{1 \leq s, t \leq n} \cos(2\pi w(t-s)w) \gamma(t-s).\end{aligned}$$

It is not difficult to check that, for any  $-(n-1) \leq h \leq n-1$ , the number of  $(s, t)$  pairs when  $t-s = h$ ,  $1 \leq s, t \leq n$ , equals  $n - |h|$ , and therefore

$$\begin{aligned}\mathrm{E}[\tilde{X}_{c,n}(w)]^2 + \mathrm{E}[\tilde{X}_{s,n}(w)]^2 &= n^{-1} \sum_{h=-n+1}^{n-1} (n - |h|) \cos(2\pi hw) \gamma(h) \\ &= \sum_{h=-n+1}^{n-1} (1 - |h|/n) \cos(2\pi hw) \gamma(h).\end{aligned}$$

Clearly, the last sum is nonnegative, and since  $\sum |\gamma(h)| < \infty$ , this sum converges to  $\sum_{h=-\infty}^{\infty} \cos(2\pi hw) \gamma(h) = f(w)$  by the Dominated Convergence Theorem. Thus  $f(w) \geq 0$  for any  $|w| \leq 1/2$ .

**Example 13.2.6.** Let  $\{X_t\}$  be iid with mean 0 and variance  $\sigma^2$ . Then  $\gamma(h) = \sigma^2$  if  $h = 0$  and  $\gamma(h) = 0$  if  $h \neq 0$ . For this case, the spectral density function is  $f(w) = \sigma^2$  for any  $w \in [-1/2, 1/2]$ .

**Example 13.2.7.** Let  $\{X_t\}$  be a mean zero MA(1) series, that is,  $X_t = \varepsilon_t + \theta\varepsilon_{t-1}$  where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ . We have already seen that

$$\gamma(h) = \begin{cases} (1 + \theta^2)\sigma^2 & h = 0 \\ \theta\sigma^2 & h = 1 \\ 0 & h \geq 2 \end{cases}$$

Since  $\gamma(h) = \gamma(-h)$  for any integer  $h$ , the spectral density function is

$$\begin{aligned} f(w) &= \gamma(0) + 2 \sum_{h=1}^{\infty} \cos(2\pi hw)\gamma(h) \\ &= (1 + \theta^2)\sigma^2 + 2\theta\sigma^2 \cos(2\pi w) \\ &= \sigma^2 \left[ 1 + \theta^2 + 2\theta \cos(2\pi w) \right]. \end{aligned}$$

**Example 13.2.8.** Let  $\{X_t\}$  be a stationary AR(1) series as described in [Example 13.2.3](#). We have already seen that for this sequence  $\gamma(h) = \phi^h\gamma(0)$ ,  $h \geq 0$ , and  $\gamma(0) = (1 - \phi^2)^{-1}\sigma^2$ . So the spectral density function of the AR(1) series is

$$\begin{aligned} f(w) &= \gamma(0) + 2 \sum_{h=1}^{\infty} \cos(2\pi hw)\gamma(h) \\ &= \gamma(0) + 2 \sum_{h=1}^{\infty} \phi^h\gamma(0) \cos(2\pi hw) \\ &= \gamma(0) \left[ 1 + 2 \sum_{h=1}^{\infty} \phi^h \cos(2\pi hw) \right]. \end{aligned}$$

Denoting  $z = \exp(-2\pi iw)$  and noting that  $2 \cos(2\pi hw) = z^h + \bar{z}^h$ , where  $\bar{z} = \exp(2\pi iw)$  is the complex conjugate of  $z$ , we have

$$\begin{aligned} 1 + 2 \sum_{h=1}^{\infty} \phi^h \cos(2\pi hw) &= 1 + \sum_{h=1}^{\infty} \phi^h (z^h + \bar{z}^h) \\ &= \sum_{h=0}^{\infty} (\phi z)^h + \sum_{h=0}^{\infty} (\phi \bar{z})^h - 1 = (1 - \phi z)^{-1} + (1 - \phi \bar{z})^{-1} - 1 \\ &= (1 - \phi^2) \frac{1}{1 + \phi^2 - 2\phi \cos(2\pi w)}, \end{aligned}$$

where the last step is obtained after some simple algebra. Hence we have

$$\begin{aligned} f(w) &= \gamma(0)(1 - \phi^2) \frac{1}{1 + \phi^2 - 2\phi \cos(2\pi w)} \\ &= \sigma^2 \frac{1}{1 + \phi^2 - 2\phi \cos(2\pi w)}, \end{aligned}$$

where the justification of the last step follows from the fact that  $\gamma(0) = (1 - \phi^2)^{-1}\sigma^2$ .

*Remark 13.2.1.* There is a more general result on the representation of the autocovariance function which does not require a restrictive condition like  $\sum|\gamma(h)| < \infty$ . It can be shown (Bochner's Theorem) that there exists a nonincreasing function  $F$  (called the spectral distribution function) on  $[-1/2, 1/2]$  such that  $\gamma(h) = \int_{-1/2}^{1/2} \exp(2\pi i h w) dF(w)$  for any integer  $h$ . When  $F$  is differentiable, its derivative  $f$  is the spectral density function. Further details can be found in the books by Brockwell and Davis [63] and Gikhman and Skorokhod [64].

### 13.2.2 Linear Time Series

In this book we only consider linear time series, which, as the same suggests, are linear combinations of mean zero variables  $\{\varepsilon_t\}$  of the form

$$X_t = \mu + \sum_{j=-\infty}^{\infty} a_j \varepsilon_{t-j},$$

where  $\sum a_j^2 < \infty$ , and  $\{\varepsilon_t\}$  are mutually uncorrelated with mean 0 and variance  $\sigma^2$ . For this series, the mean and the autocovariance function are

$$\mathbb{E}[X_t] = \mu, \quad \gamma(h) = \sigma^2 \sum_{j=-\infty}^{\infty} a_j a_{j+h}.$$

Such a representation exists if the spectral density  $f$  is integrable. If  $\{X_t\}$  is stationary Gaussian, then  $\{\varepsilon_t\}$  are iid  $N(0, \sigma^2)$  variables.

A stationary series  $\{X_t\}$  is said to have a causal representation if it can be written as

$$X_t = \mu + \sum_{j=0}^{\infty} a_j \varepsilon_{t-j},$$

where  $\{\varepsilon_t\}$  are mutually uncorrelated mean zero rv's with common variance  $\sigma^2$ . If  $\int_{-1/2}^{1/2} \log(f(w)) dw > -\infty$ , then the series admits a causal representation. As before, if  $\{X_t\}$  is stationary Gaussian, then  $\{\varepsilon_t\}$  are iid  $N(0, \sigma^2)$  variables. A detailed discussion on these representations can be found in [Chapter 5](#) of the book by Gikhman and Skorokhod [64].

In practice, it is usually assumed that  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$  even if they are not normally distributed. Particular examples of linear time series are ARMA processes. There are other time series that are nonlinear such as the bilinear series

$$X_t = \mu + \sum_{j=-\infty}^{\infty} a_j \varepsilon_{t-j} + \sum_{-\infty < j, k < \infty} b_{jk} \varepsilon_{t-j} \varepsilon_{t-k}$$

with  $\sum |a_j| < \infty$  and  $\sum |b_{jk}| < \infty$ , but such series are not discussed in this chapter.

From now on, it will be assumed that  $\{X_t\}$  is a linear time series with  $\{\varepsilon_t\}$  iid mean zero rv's with finite variance  $\sigma^2$ , that is,

$$X_t = \mu + \sum_{j=-\infty}^{\infty} a_j \varepsilon_{t-j}, \quad \text{with } \sum_{j=-\infty}^{\infty} |a_j| < \infty, \quad (2)$$

$\{\varepsilon_t\}$  iid,  $E[\varepsilon_t] = 0$ , and  $\text{Var}[\varepsilon_t] = \sigma^2$ .

This series  $\{X_t\}$  is stationary as discussed above.

### Notations

The following notations will be used throughout this chapter. For any stationary series with autocovariance function  $\gamma$  and for any positive integer  $h$ ,  $\boldsymbol{\Gamma}_h$  will denote the  $h \times h$  matrix whose element  $(j, k)$  is given by  $\gamma(j - k)$  and  $\boldsymbol{\gamma}_h$  will denote the  $h \times 1$  vector with elements  $\gamma(1), \dots, \gamma(h)$ . Similarly,  $\hat{\boldsymbol{\Gamma}}_h$  and  $\hat{\boldsymbol{\gamma}}_h$  will denote the estimates of  $\boldsymbol{\Gamma}_h$  and  $\boldsymbol{\gamma}_h$  when the autocovariances are estimated based on the available data. We note that  $\boldsymbol{\Gamma}_h$  is the covariance matrix of  $(X_{t+1}, \dots, X_{t+h})$  and it is nonnegative definite.

#### 13.2.3 Time Reversibility for Linear Prediction

If  $\{X_t\}$  is stationary with autocovariance function  $\gamma$ , then we show below that the coefficients associated with the best linear predictor (forecast) of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$  are the same as those of the best linear predictor (backcast) of  $X_{t-h-1}$  based on  $X_{t-h}, \dots, X_{t-1}$ . This result is true for any  $h \geq 1$ . In order to simplify the notations, we assume that  $\{X_t\}$  has zero mean, since we can always subtract the mean  $\mu$  from  $X_t$  to achieve this.

A linear predictor  $\phi_1 X_{t-1} + \dots + \phi_h X_{t-h}$  of  $X_t$  has the prediction error  $PE^{(f)} = E[X_t - \phi_1 X_{t-1} - \dots - \phi_h X_{t-h}]^2$ . Since  $X_t, X_{t-1}, \dots, X_{t-h}$  have zero means, there is no need to include an intercept term in the formula for the predictor. The best linear predictor can be obtained by minimizing  $PE^{(f)}$  with respect to  $\phi_1, \dots, \phi_h$ . As in the case of linear models, we can differentiate  $PE^{(f)}$  with respect to  $\phi_j, j = 1, \dots, h$ , and equate the derivatives to zero, which leads to the normal equations. When the derivative of  $PE^{(f)}$  with respect to  $\phi_j$  is set to 0, we have

$$\begin{aligned} -2E[X_t - \phi_1 X_{t-1} - \dots - \phi_h X_{t-h}]X_{t-j} &= 0, \text{ ie,} \\ -\gamma(j) + \phi_1 \gamma(j-1) + \dots + \phi_h \gamma(j-h) &= 0, \text{ ie,} \\ \phi_1 \gamma(j-1) + \dots + \phi_h \gamma(j-h) &= \gamma(j), \end{aligned}$$

$j = 1, \dots, h$ . These normal equations can be written in the matrix form as

$$\boldsymbol{\Gamma}_h \boldsymbol{\phi} = \boldsymbol{\gamma}_h, \quad (3)$$

where  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_h)^T$ , and they are also known as the *Yule-Walker equations*. Moreover, the prediction error of the best linear predictor with coefficients  $\phi_1, \dots, \phi_h$  satisfying the Yule-Walker equations is given by

$$\begin{aligned}
E[X_t - \phi_1 X_{t-1} - \cdots - \phi_h X_{t-h}]^2 &= E[X_t^2] - 2E[X_t(\phi_1 X_{t-1} + \cdots + \phi_h X_{t-h})] \\
&\quad + E[\phi_1 X_{t-1} + \cdots + \phi_h X_{t-h}]^2 \\
&= \gamma(0) - 2\boldsymbol{\phi}^T \boldsymbol{\gamma}_h + \boldsymbol{\phi}^T \boldsymbol{\Gamma}_h \boldsymbol{\phi} = \gamma(0) - \boldsymbol{\phi}^T \boldsymbol{\gamma}_h \\
&= \gamma(0) - \phi_1 \gamma(1) - \cdots - \phi_h \gamma(h).
\end{aligned}$$

A linear predictor  $\phi_1 X_{t-h} + \phi_2 X_{t-h+1} + \cdots + \phi_h X_{t-1}$  of  $X_{t-h-1}$  has the prediction error  $PE^{(b)} = E[X_{t-h-1} - \phi_1 X_{t-h} - \cdots - \phi_h X_{t-1}]^2$ . In order to obtain the best linear predictor we need to differentiate  $PE^{(b)}$  with respect to  $\phi_1, \dots, \phi_h$  and equate the derivatives to 0. Setting the derivative of  $PE^{(b)}$  with respect to  $\phi_j$  to 0, we have

$$\begin{aligned}
-2E[(X_{t-h-1} - \phi_1 X_{t-h} - \cdots - \phi_h X_{t-1}) X_{t-h-1+j}] &= 0, \text{ ie,} \\
-\gamma(j) + \phi_1 \gamma(j-1) + \cdots + \phi_h \gamma(j-h) &= 0, \text{ ie,} \\
\phi_1 \gamma(j-1) + \cdots + \phi_h \gamma(j-h) &= \gamma(j),
\end{aligned}$$

$j = 1, \dots, h$ . These normal equations in the matrix form are exactly the same as in Eq. (3). Since the normal equations for forecasting (ie, linear prediction of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ ) and backcasting (ie, linear prediction of  $X_{t-h-1}$  based on  $X_{t-h}, \dots, X_{t-1}$ ) are the same, the vectors of coefficients  $\boldsymbol{\Gamma}_h^{-1} \boldsymbol{\gamma}_h$  are also the same. Using the arguments given above, it is fairly easy to verify that the prediction error for the best linear predictor for backcasting  $X_{t-h-1}$  is exactly the same as that for forecasting  $X_t$ . Thus we arrive at an important result.

**Lemma 13.2.1.** *Let  $\{X_t\}$  be a mean zero stationary series.*

- (a) *For any positive integer  $h$ , if  $\hat{X}_t^{(f)} = \phi_1 X_{t-1} + \cdots + \phi_h X_{t-h}$  is the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ , then  $\hat{X}_{t-h-1}^{(b)} = \phi_1 X_{t-h} + \phi_2 X_{t-h+1} + \cdots + \phi_h X_{t-1}$  is the best linear predictor of  $X_{t-h-1}$  based on  $X_{t-h}, \dots, X_{t-1}$ .*
- (b) *The vector of coefficients  $\boldsymbol{\phi}$  of the best predictor described in part (a) is given by the solution of the normal equations (ie, the Yule-Walker equations)  $\boldsymbol{\Gamma}_h \boldsymbol{\phi} = \boldsymbol{\gamma}_h$ .*
- (c) *The prediction errors  $E[X_t - \hat{X}_t^{(f)}]^2$  and  $E[X_{t-h-1} - \hat{X}_{t-h-1}^{(b)}]^2$  are the same and they are given by  $\gamma(0) - \phi_1 \gamma(1) - \cdots - \phi_h \gamma(h)$ .*

### 13.3 Estimation of the Mean and the Autocorrelation Function

Suppose that  $X_1, \dots, X_n$  are observations from a stationary series with mean  $\mu$  and covariance function  $\gamma$ . Some useful descriptive statistics for an initial analysis involve estimation of the mean, the autocovariance function  $\gamma$ , the autocorrelation function  $\rho$ , etc. These are the topics of discussion in this section.

#### 13.3.1 Estimation of the Mean

An estimate of  $\mu$  is given by  $\hat{\mu} = \bar{X}_n = n^{-1} \sum_{t=1}^n X_t$ . It is fairly easy to see that  $\hat{\mu}$  is an unbiased estimator of  $\mu$ . In order to obtain the standard error of this estimate, one needs to calculate the variance of  $\hat{\mu}$ . Note that

$$\text{Var}[\hat{\mu}] = n^{-2} \sum_{s,t=1}^n \text{Cov}[X_s, X_t] = n^{-2} \sum_{s,t=1}^n \gamma(s-t).$$

The argument used in [Section 13.2.1](#) to demonstrate the nonnegativeness of the spectral density function can be employed to show that

$$\begin{aligned} \tau_n^2 &= n\text{Var}[\hat{\mu}] = n^{-1} \sum_{h=-(n-1)}^{n-1} (n - |h|)\gamma(h) = \sum_{h=-n+1}^{n-1} (1 - |h|/n)\gamma(h) \\ &\rightarrow \sum_{h=-\infty}^{\infty} \gamma(h) := \tau^2, \end{aligned}$$

assuming that  $\sum |\gamma(h)| < \infty$ . Since  $\gamma(-h) = \gamma(h)$ , we may also rewrite  $\tau_n^2$  and  $\tau^2$  as

$$\begin{aligned} \tau_n^2 &= n\text{Var}[\hat{\mu}] = \sum_{h=-n+1}^{n-1} (1 - |h|/n)\gamma(h) = \gamma(0) + 2 \sum_{h=1}^{n-1} (1 - h/n)\gamma(h), \\ \tau^2 &= \sum_{h=-\infty}^{\infty} \gamma(h) = \gamma(0) + 2 \sum_{h=1}^{\infty} \gamma(h). \end{aligned}$$

If it is assumed that  $\gamma(h)$  converges rapidly to zero as  $h \rightarrow \infty$ , a condition satisfied by many series such as  $ARMA(p, q)$ , then one may ignore  $\gamma(h)$ ,  $h > L$ , for a suitably chosen integer  $L$  which may depend on  $n$  (eg,  $L \approx \sqrt{n}$ ). In such a case, a reasonable estimate of  $\tau_n^2$  is

$$\hat{\tau}_n^2 = n^{-1} \hat{\gamma}(0) + 2 \sum_{h=1}^L (1 - h/n) \hat{\gamma}(h) \approx n^{-1} \left[ \hat{\gamma}(0) + 2 \sum_{h=1}^L \hat{\gamma}(h) \right],$$

where  $\hat{\gamma}(h)$  is an estimate of  $\gamma(h)$ . It turns out that the central limit theorem (CLT) holds for  $\hat{\mu} = \bar{X}_n$  under reasonable conditions.

**Theorem 13.3.1.** *Assume that  $X_1, \dots, X_n$  are observations from a linear stationary series  $\{X_t\}$  as given in Eq. (2). Then as  $n \rightarrow \infty$ ,  $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{\mathcal{D}} N(0, \tau^2)$ , where  $\tau^2 = \lim_{n \rightarrow \infty} n\text{Var}[\hat{\mu}]$ .*

This result now can be used to construct an approximate confidence interval for  $\mu$ . Let  $\hat{\tau}_n$  be as given above. Assuming that  $\hat{\tau}_n/\tau_n - 1 \xrightarrow{P} 0$ , an approximate confidence interval for  $\mu$  with confidence coefficient  $1 - \alpha$  is  $\hat{\mu} \pm z_{\alpha/2} \hat{\tau}_n/\sqrt{n}$ .

A proof of the above result on the asymptotic normality of  $\hat{\mu} = \bar{X}_n$  uses the CLT for the sample mean of  $m$ -dependent rv's ([Theorem 3.3.3](#)) and it involves careful details as it requires verifying technical conditions, and the details are given in Brockwell and Davis [63]. However, the basic idea behind the proof is simple and can be summarized in the following steps. Let  $W_n = \sqrt{n}(\bar{X}_n - \mu)$ .

I. The truncated rv's  $\{X_{m,t}\}$

$$X_{m,t} = \mu + \sum_{j=-m}^m a_j \varepsilon_{t-j},$$

are  $2m$ -dependent with mean  $\mu$  where  $m$  is a positive integer.

- II. Let  $\bar{X}_{m,n} = n^{-1} \sum_{t=1}^n X_{m,t}$  and  $W_{m,n} = \sqrt{n}(\bar{X}_{m,n} - \mu)$ . Since  $\{X_{m,t}\}$  is  $2m$ -dependent, the CLT holds for  $W_{m,n}$  (Theorem 3.3.3). Thus  $W_{m,n} \xrightarrow{D} Z_m \sim N(0, \tau(m)^2)$ , where  $\tau^2(m) = \lim_{n \rightarrow \infty} \text{Var}[W_{m,n}]$ .
- III. Show that for any  $\delta > 0$ ,

$$\limsup_{n \rightarrow \infty} P[|W_{m,n} - W_n| > \delta] \rightarrow 0$$

as  $m \rightarrow \infty$ . A sufficient condition for this, via Chebychev's inequality, is

$$\limsup_{n \rightarrow \infty} E[|W_{m,n} - W_n|^2] \rightarrow 0$$

as  $m \rightarrow \infty$ .

- IV. If  $\tau^2(m) = \lim_{n \rightarrow \infty} \text{Var}[W_{m,n}] \rightarrow \tau^2 > 0$ , then  $W_n \xrightarrow{D} Z \sim N(0, \tau^2)$ .

The four steps given above can be written down in a general framework which does not require  $\{X_t\}$  to be a linear stationary series. The general framework is useful since this result can also be applied for obtaining asymptotic normality of estimates of the autocovariances and autocorrelations.

**Theorem 13.3.2.** *Let  $X_1, \dots, X_n$  be observations from an infinite sequence  $\{X_t\}$  of rv's with common mean  $\mu$  and let  $W_n = \sqrt{n}(\bar{X}_n - \mu)$ , where  $\bar{X}_n = n^{-1} \sum_{t=1}^n X_t$ . Assume that for any positive integer  $m$ , there exists an  $m$ -dependent series  $\{X_{m,t}\}$  with common mean  $\mu_m$  and let  $W_{m,n} = \sqrt{n}(\bar{X}_{m,n} - \mu_m)$ , where  $\bar{X}_{m,n} = n^{-1} \sum_{t=1}^n X_{m,t}$ . Then  $W_n \xrightarrow{D} Z \sim N(0, \tau^2)$  if the following conditions hold*

- (i) for every positive integer  $m$ ,  $W_{m,n} \xrightarrow{D} Z_m \sim N(0, \tau(m)^2)$ , where  $\tau^2(m) = \lim_{n \rightarrow \infty} \text{Var}[W_{m,n}]$ ,
- (ii) for any  $\delta > 0$ ,  $\limsup_{n \rightarrow \infty} P[|W_{m,n} - W_n| > \delta] \rightarrow 0$  as  $m \rightarrow \infty$ ,
- (iii)  $\tau^2(m)$  converges to a constant  $\tau^2 > 0$  as  $m \rightarrow \infty$ .

Let us briefly examine why this theorem is true. For any real number  $z$  and  $\delta > 0$ ,

$$\begin{aligned} P[W_n \leq z] &= P[W_n \leq z, |W_n - W_{m,n}| \leq \delta] \\ &\quad + P[W_n \leq z, |W_n - W_{m,n}| > \delta] \\ &\leq P[W_{m,n} \leq z + \delta] + P[|W_n - W_{m,n}| > \delta]. \end{aligned}$$

Denoting  $\Delta(m, \delta) = \limsup_{n \rightarrow \infty} P[|W_n - W_{m,n}| > \delta]$ , we have

$$\begin{aligned} \limsup_{n \rightarrow \infty} P[W_n \leq z] &\leq \limsup_{n \rightarrow \infty} P[W_{m,n} \leq z + \delta] + \Delta(m, \delta) \\ &= P[Z_m \leq z + \delta] + \Delta(m, \delta). \end{aligned}$$

Since  $\tau(m)^2 \rightarrow \tau^2$ , letting  $m \rightarrow \infty$  we have

$$\begin{aligned}\limsup_{n \rightarrow \infty} P[W_n \leq z] &\leq P[Z \leq z + \delta], \text{ and thus} \\ \limsup_{n \rightarrow \infty} P[W_n \leq z] &\leq P[Z \leq z],\end{aligned}$$

since  $\delta > 0$  is arbitrary.

A similar argument will show that

$$\begin{aligned}\limsup_{n \rightarrow \infty} P[W_n > z] &\leq P[Z > z], \text{ and hence} \\ \liminf_{n \rightarrow \infty} P[W_n \leq z] &\geq P[Z \leq z].\end{aligned}$$

Since

$$P[Z \leq z] \leq \liminf_{n \rightarrow \infty} P[W_n \leq z] \leq \limsup_{n \rightarrow \infty} P[W_n \leq z] \leq P[Z \leq z],$$

we conclude that  $\lim_{n \rightarrow \infty} P[W_n \leq z] = P[Z \leq z]$ .

### 13.3.2 Estimation of Autocovariance and Autocorrelation Functions

Estimates of  $\gamma(h) = \text{Cov}[X_t, X_{t+h}]$  and  $\rho(h) = \text{Corr}[X_t, X_{t+h}]$ ,  $h = 0, 1, \dots$  are

$$\hat{\gamma}(h) = n^{-1} \sum_{t=1}^{n-h} (X_t - \bar{X})(X_{t+h} - \bar{X}) \quad \text{and} \quad \hat{\rho}(h) = \hat{\gamma}(h)/\hat{\gamma}(0). \quad (4)$$

This estimate of  $\gamma(h)$  is not unbiased as will be clear soon. For notational convenience, we assume that  $E[X_t] = 0$  since  $X_t - \bar{X}$  is the same as  $Y_t - \bar{Y}$ , where  $Y_t = X_t - \mu$  is the centered  $X_t$ . Simple algebra shows that

$$\begin{aligned}\hat{\gamma}(h) &= n^{-1} \sum_{t=1}^{n-h} (X_t - \bar{X})(X_{t+h} - \bar{X}) \\ &= n^{-1} \sum_{t=1}^{n-h} X_t X_{t+h} - \bar{X} n^{-1} \sum_{t=1}^{n-h} X_{t+h} - \bar{X} n^{-1} \sum_{t=1}^{n-h} X_t + n^{-1} (n-h) \bar{X}^2 \\ &= n^{-1} \sum_{t=1}^{n-h} X_t X_{t+h} - (1 + h/n) \bar{X}^2 + \bar{X} \left[ n^{-1} \sum_{t=1}^h X_t + n^{-1} \sum_{t=n-h+1}^n X_t \right].\end{aligned}$$

We have already seen from the last section that  $\text{Var}[\bar{X}] = \tau_n^2/n = O(n^{-1})$  assuming that  $\sum |\gamma(h)| < \infty$ . Since  $E[X_t] = 0$ ,  $E[\bar{X}^2] = \text{Var}[\bar{X}]$ . It is not difficult to verify that the expectation of the last term in the last displayed equation is  $O(n^{-3/2})$ . Since  $E[X_t X_{t+h}] = \gamma(h)$ , we have

$$\begin{aligned} E[\hat{\gamma}(h)] &= n^{-1}(n-h)\gamma(h) + \tau_n^2/n + O(n^{-3/2}) \\ &= \gamma(h) - (h/n)\gamma(h) + \tau_n^2/n + O(n^{-3/2}). \end{aligned}$$

Expression for  $\tau_n^2/n$  given in the last section shows that  $(h/n)\gamma(h)$  is not cancelled by  $\tau_n^2/n$  and thus the bias of  $\hat{\gamma}(h)$  is  $O(h/n)$ .

From the preceding arguments, we get the following simple result.

**Lemma 13.3.1.** *Let  $\{X_t\}$  be stationary with the autocovariance function  $\{\gamma(h)\}$  which satisfies the condition  $\sum |\gamma(h)| < \infty$ . Then*

$$\begin{aligned} \hat{\gamma}(h) &= n^{-1} \sum_{t=1}^{n-h} (X_t - \mu)(X_{t+h} - \mu) - (\bar{X} - \mu)^2 + R_1 \\ &= n^{-1} \sum_{t=1}^n (X_t - \mu)(X_{t+h} - \mu) + R_2, \end{aligned}$$

where  $E|R_1| = O(n^{-3/2})$  and  $E|R_2| = O(n^{-1})$ .

Why is the biased estimate of  $\gamma(h)$  used in practice? The reason is that the  $n \times n$  matrix  $\hat{\Gamma}_n = ((\hat{\gamma}(j-k)))_{n \times n}$ , which is an estimate of the covariance matrix  $\Gamma_n$ , is nonnegative definite. Had we used an unbiased estimate of  $\gamma(h)$ , then this nonnegative definiteness property is not guaranteed to be preserved. An unbiased estimate of  $\gamma(h)$  is

$$\begin{aligned} \tilde{\gamma}(h) &= (n-h-1)^{-1} \sum_{t=1}^{n-h} (X_t - \bar{X}_1)(X_{t+h} - \bar{X}_2), \text{ with} \\ \bar{X}_1 &= (n-h)^{-1} \sum_{t=1}^{n-h} X_t \text{ and } \bar{X}_2 = (n-h)^{-1} \sum_{t=h+1}^n X_t. \end{aligned}$$

Asymptotic distributions of the estimated autocovariances and autocorrelations are known. We first write down the joint distribution of  $\hat{\rho}(1), \dots, \hat{\rho}(h)$ .

**Theorem 13.3.3.** *Let  $\{X_t\}$  be a linear stationary series as given in Eq. (2) with the extra assumption  $E[\varepsilon_t^4] < \infty$ . Let  $\rho_h$  be the  $h \times 1$  vector with elements  $\rho(1), \dots, \rho(h)$ , and similarly let  $\hat{\rho}_h$  be the  $h \times 1$  vector with elements  $\hat{\rho}(1), \dots, \hat{\rho}(h)$ , where  $\hat{\rho}(j)$ 's are as in Eq. (4). Then  $\sqrt{n}(\hat{\rho} - \rho) \xrightarrow{\mathcal{D}} N_h(\mathbf{0}, \mathbf{W})$  as  $n \rightarrow \infty$ , where element  $(j, k)$  of the matrix  $\mathbf{W}$  is*

$$\begin{aligned} w_{jk} &= \sum_{l=-\infty}^{\infty} e(j, l)e(k, l), \text{ with} \\ e(j, l) &= [\rho(l+j) + \rho(l-j) - 2\rho(j)\rho(l)]/\sqrt{2}. \end{aligned}$$

If  $\{X_t\}$  are iid, then  $\rho(j) = 0$  for any  $j \geq 1$ ,  $\mathbf{W}$  is the  $h \times h$  identity matrix, and,  $\sqrt{n}\hat{\rho}(j)$ ,  $j = 1, \dots, h$ , are asymptotically iid  $N(0, 1)$  variables.

Note that if the series  $\{X_t\}$  are iid, then  $e(j, l)$  is equal to 1 if  $l = j$  and is 0 otherwise. In such a case, the matrix  $\mathbf{W}$  in the above theorem is the identity. So if  $\{X_t\}$  are iid, then  $\sqrt{n}\hat{\rho}(j)$ ,  $j = 1, 2, \dots$  are asymptotically iid  $N(0, 1)$  rv's. The estimated autocorrelation function is widely used to check if a series is white noise, that is, to check if  $\rho(h) = 0$  when  $h \neq 0$ . In

practice, one usually plots  $\hat{\rho}(h)$  against  $h$  to check if  $\hat{\rho}(h)$  is outside the range  $0 \pm 2/\sqrt{n}$  to investigate if  $\rho(h) \neq 0$ . Such a plot is called the ACF plot.

The proof of the above theorem on the asymptotic distribution of the estimated autocorrelations follows from the asymptotic normality of the estimated autocovariances and this is stated below.

**Theorem 13.3.4.** *Let  $\{X_t\}$  be a linear stationary series as given in Eq. (2) with the extra assumption  $E[\varepsilon_t^4] < \infty$ . Let  $\delta_h$  be the  $(h+1)$ -dim vector with elements  $\gamma(0), \dots, \gamma(h)$  and, similarly, let  $\hat{\delta}_h$  be the vector of  $\hat{\gamma}(0), \dots, \hat{\gamma}(h)$ . Then  $\sqrt{n}(\hat{\delta}_h - \delta_h) \xrightarrow{\mathcal{D}} N_{h+1}(\mathbf{0}, \mathbf{V})$  where element  $(j, k)$  of  $\mathbf{V}$  is*

$$\begin{aligned} & [E(\varepsilon_t^4) - 3\sigma^4] \gamma(j)\gamma(k) + \sum_{l=-\infty}^{\infty} f(j, l)f(k, l), \text{ with} \\ & \sigma^2 = E[\varepsilon_t^2], f(j, l) = [\gamma(l+j) + \gamma(l-j)]/\sqrt{2}. \end{aligned}$$

The proof of the above result uses the following steps and [Theorem 13.3.2](#) given above. The details are quite long, and are given in the book by Brockwell and Davis [63].

As in the analysis of bias of  $\hat{\gamma}(h)$ , we may use the centered variable  $X_t - \mu$ , but we continue to denote it by  $X_t$  with the understanding that  $E[X_t] = 0$ .

- I. Define  $\tilde{\gamma}(j) = n^{-1} \sum_{t=1}^n X_t X_{t+j}$ . Then, by [Lemma 13.3.1](#),  $\hat{\gamma}(j) = \tilde{\gamma}(j) + O_P(n^{-1})$  since  $\bar{X} = O_P(n^{-1/2})$ .
- II. It is enough to prove the CLT for  $\sum_{j=0}^h c_j \hat{\gamma}(j)$  where  $c_0, \dots, c_h$  are constants (Cramér-Wold device). Since  $\hat{\gamma}(j) = \tilde{\gamma}(j) + O_P(n^{-1})$ ,  $j = 0, \dots, h$ , it is enough to prove the CLT for  $\sum_{j=0}^h c_j \tilde{\gamma}(j)$ .
- III. Denoting  $S_t = X_t(c_0 X_t + \dots + c_h X_{t+h})$ , we have

$$\sum_{j=0}^h c_j \tilde{\gamma}(j) = n^{-1} \sum_{t=1}^n X_t(c_0 X_t + \dots + c_h X_{t+h}) = n^{-1} \sum_{t=1}^n S_t.$$

- IV. Let  $X_{m,t} = \sum_{j=-m}^m a_j \varepsilon_{t-j}$ , where  $m$  is a positive integer and

$$S_{m,t} = X_{m,t}(c_0 X_{m,t} + \dots + c_h X_{m,t+h}).$$

Note that  $\{X_{m,t}\}$  and  $\{S_{m,t}\}$  are  $2m$ - and  $(h+2m)$ -dependent sequences of rv's, respectively. Moreover,  $E[S_t]$  is the same for all  $t$  and the same is true for  $E[S_{m,t}]$ . Let

$$\bar{S}_n = n^{-1} \sum_{t=1}^n S_t, \quad \bar{S}_{m,n} = n^{-1} \sum_{t=1}^n S_{m,t},$$

$$\theta = E[S_t], \quad \theta_m = E[S_{m,t}],$$

$$W_n = \sqrt{n}(\bar{S}_n - \theta), \quad \text{and} \quad W_{m,n} = \sqrt{n}(\bar{S}_{m,n} - \theta_m).$$

- V. For any positive integer  $m$ , show that  $W_{m,n} \xrightarrow{\mathcal{D}} Z_m \sim N(0, v_m^2)$  as  $n \rightarrow \infty$  (Theorem 3.4.3).
- VI. Show that  $\lim_{m \rightarrow \infty} \Delta(m, \delta) = 0$ , where  $\Delta(m, \delta) = \limsup_{n \rightarrow \infty} P[|W_n - W_{m,n}| > \delta]$ .
- VII. Show that  $v_m \rightarrow v^2$  as  $m \rightarrow \infty$ , where  $v^2 = \sum_{0 \leq j, k \leq h} c_j c_k v_{jk}$  ( $v_{jk}$  is element  $(j, k)$  of the matrix  $V$ ).

### 13.3.3 Diagnostics

Diagnostic tools are often employed for an observed series  $X_1, \dots, X_n$  in order to investigate if there is trend, if the assumption of equal variance is reasonable, if the sequence is iid, if the assumption of normality is reasonable, etc. Before we discuss the diagnostic methods, let us introduce another descriptive measure called the partial autocorrelation function (PACF) which is commonly used in the analysis of stationary time series data. PACF of order  $h$  is the partial correlation between  $X_t$  and  $X_{t-h}$  given  $X_{t-1}, \dots, X_{t-h+1}$ . More formally, if  $X_t^{(f)}$  is the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h+1}$  and  $X_{t-h}^{(b)}$  is the best linear predictor of  $X_{t-h}$  based on  $X_{t-h+1}, \dots, X_{t-1}$ , then the partial autocorrelation between  $X_t$  and  $X_{t-h}$  is  $\pi(h) = \text{Corr}[X_t - X_t^{(f)}, X_{t-h} - X_{t-h}^{(b)}]$ ,  $h \geq 2$ . By definition,  $\pi(1) = \rho(1)$ . In a later section, partial autocorrelations will be discussed in detail along with appropriate formulas that can be used for computations. If a series is  $AR(p)$ , then  $\pi(h) = 0$  when  $h \geq p + 1$ , and if  $\hat{\pi}(h)$  is the estimator of  $\pi(h)$ , then  $\sqrt{n}\hat{\pi}(h)$  is approximately distributed as  $N(0, 1)$  for large  $n$ .

We now discuss some graphical and formal diagnostic procedures.

- (a) Plot of the series against time reveals if there is a trend or if the assumption of equal variance (across time) is reasonable. In some cases, the problem of unequal variance can be remedied by an appropriate Box-Cox transformation of the observed series.
- (b) In order to check if  $\{X_t\}$  are iid, one may plot the estimated autocorrelations (the ACF plot) along with  $0 \pm 2/\sqrt{n}$  bars in order to assess if the autocorrelations are close to zero. ACF plot of the estimated residuals  $\{\hat{\varepsilon}_t\}$  after fitting an ARMA model can also be used to assess appropriateness of the model. For instance, if an  $AR(p)$  model is fitted, the estimated residuals are

$$\hat{\varepsilon}_t = X_t - \bar{X} - \sum_{j=1}^p \hat{\phi}_j (X_{t-j} - \bar{X}), \quad t = p+1, \dots, n,$$

where  $\bar{X}$  is the sample mean, and  $\hat{\phi}_1, \dots, \hat{\phi}_p$  are the estimated autoregressive parameters. If  $AR(p)$  provides an appropriate description of the data, then  $\{\hat{\varepsilon}_t\}$  are approximately iid and its ACF plot would indicate this. Similar logic can be used when fitting an  $MA(q)$  or an  $ARMA(p, q)$  model to check its adequacy.

- (c) There are a number of formal tests for checking if the observations are iid. We may want to test, for a given positive integer  $h$ ,  $H_0: \rho(1) = \dots = \rho(h) = 0$  against  $H_1:$  at

least one of  $\rho(1), \dots, \rho(h)$  is nonzero. We mention two tests: Portmanteau and Ljung-Box. The test statistics are

$$Q = n \sum_{j=1}^h \hat{\rho}(j)^2 \quad \text{and} \quad Q_{LB} = n(n+2) \sum_{j=1}^h \hat{\rho}(j)^2 / (n-j),$$

where  $\{\hat{\rho}(j)\}$  are the estimated autocorrelations as given in Eq. (4). Under  $H_0$ , each of the two statistics ( $Q$  and  $Q_{LB}$ ) has an approximate chi-squared distribution with  $h$  degrees of freedom. So we can reject the null hypothesis at level  $\alpha$  if  $Q > \chi_{h,\alpha}^2$  or  $Q_{LB} > \chi_{h,\alpha}^2$ , where  $\chi_{h,\alpha}^2$  is the  $(1 - \alpha)$ -quantile of the chi-squared distribution with  $h$  degrees of freedom. These tests are based on the asymptotic result given in **Theorem 13.3.3**. According to this theorem, if the series consists of iid observations, then  $\sqrt{n}\hat{\rho}(j)$ ,  $j = 1, \dots, h$ , are approximately iid  $N(0, 1)$  for large  $n$ . Ljung-Box test provides a small sample correction to the Portmanteau test by obtaining a better estimate of the asymptotic variance of  $\hat{\rho}(j)$ . In Portmanteau test,  $\text{Var}[\hat{\rho}(j)]$  is approximated by  $1/n$ , whereas Ljung-Box uses the approximation  $(n-j)/[n(n+2)]$ .

- (d) One can examine the histogram of the data to check if the assumption of normality is justifiable.
- (e) A plot of  $\hat{\pi}(h)$  against  $h$  along with  $0 \pm 2/\sqrt{n}$  bars is known as the PACF plot and it can be used to make an assessment if  $\pi(h)$  is substantially different from zero. This plot is useful in guessing the order of an autoregressive model since  $\pi(h) = 0$ ,  $h \geq p+1$ , for an  $AR(p)$  model. It is a common practice to use ACF and PACF plots in the initial analysis of the data.

It is useful to keep in mind that for an  $MA(q)$  model, autocorrelations of lag  $q+1$  or higher are all zero. Similarly, for an  $AR(p)$  model, partial autocorrelations of order  $p+1$  or higher are zero. So the ACF plot is useful for an initial guess of the order of a moving average model and the PACF plot is useful in guessing the order of an autoregressive model. For instance, if the autocorrelations of lag 3 or higher are all negligible, then an  $MA(2)$  model may provide a reasonable description of the data. If the partial autocorrelations of lag 3 or higher are negligible, then  $AR(2)$  may be a reasonable model for the data. It should be pointed, however, out that if the true model is  $MA(q)$ , the asymptotic mean of  $\hat{\rho}(h)$  is 0 for  $h \geq q+1$ , but the asymptotic variance of  $\hat{\rho}(h)$  is  $\sum_{j=-q}^q \rho(j)^2/n$  and not  $1/n$  and thus the  $\pm 2/\sqrt{n}$  bounds are not necessarily equal to  $\pm 2SE[\hat{\rho}(h)]$ . Nevertheless, the ACF plot can be a useful graphical method for assessing if some moving average model is reasonable.

The ACF of a stationary autoregressive series, and under the condition of invertibility (discussed later), the PACF of a moving average series decrease rapidly with lag  $h$ . In general, for an ARMA series both the ACF and PACF decrease rapidly. The following table provides a summary.

Model	<i>AR(p)</i>	<i>MA(q)</i>	<i>ARMA(p, q)</i>
ACF	Tails off	Cuts off after lag $q$	Tails off
PACF	Cuts off after lag $p$	Tails off	Tails off

In order to get a good predictive model, one should use a criterion such as AIC to select an appropriate model. Nevertheless, the use of ACF and PACF plots may sometimes lead to a reasonable predictive model.

### 13.3.4 Notation of Backshift Operator

The notation of backshift operator makes the description of ARMA processes convenient. The backshift operator  $B$  is defined as  $BX_t = X_{t-1}$ . Note that  $X_{t-k} = B^k X_t$  for any  $k$  with the understanding that  $B^0 X_t = X_t$ . A mean zero *AR(1)* process can be written as

$$X_t = \phi_1 BX_t + \varepsilon_t, \text{ or } X_t - \phi_1 BX_t = \varepsilon_t, \text{ or } (1 - \phi_1 B)X_t = \varepsilon_t.$$

Similarly, a mean zero *AR(p)* process can be written as

$$(1 - \phi_1 B - \cdots - \phi_p B^p)X_t = \varepsilon_t, \text{ or } \phi(B)X_t = \varepsilon_t, \\ \text{with } \phi(B) = 1 - \phi_1 B - \cdots - \phi_p B^p.$$

If a series  $\{X_t\}$  with mean  $\mu$  is *AR(p)*, then

$$\phi(B)(X_t - \mu) = \varepsilon_t, \text{ where } \phi(B) = 1 - \phi_1 B - \cdots - \phi_p B^p.$$

Similarly, an *MA(q)* series  $\{X_t\}$  with mean  $\mu$  can be expressed as

$$X_t - \mu = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \cdots + \theta_q \varepsilon_{t-q} = \theta(B)\varepsilon_t, \text{ where} \\ \theta(B) = 1 + \theta_1 B + \cdots + \theta_q B^q.$$

An *ARMA(p, q)* series  $\{X_t\}$  with mean  $\mu$  can be written as

$$\phi(B)(X_t - \mu) = \theta(B)\varepsilon_t.$$

As mentioned earlier, the error terms  $\{\varepsilon_t\}$  are often called *innovations* in the time series literature.

## 13.4 Partial Autocorrelation Function (PACF)

As discussed in the last section, the partial autocorrelation (PACF) plot is widely used in the analysis of stationary data as it aids in the preliminary identification of autoregressive models. Let us recall that, for any  $h \geq 2$ , partial autocorrelation  $\pi(h)$  of order  $h$  of a stationary series  $\{X_t\}$  is defined to be the partial correlation between  $X_t$  and  $X_{t-h}$  given  $X_{t-1}, \dots, X_{t-h+1}$ . When  $h = 1$ ,  $\pi(1)$  is defined to be  $\text{Corr}[X_t, X_{t-1}]$ . If  $\hat{X}_{h-1,t}^{(f)}$  is the best linear predictor (forecast) of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h+1}$  and  $\hat{X}_{h-1,t-h}^{(b)}$  is the best linear predictor

(backcast) of  $X_{t-h}$  based on  $X_{t-h+1}, \dots, X_{t-1}$ , then for any  $h \geq 2$ , the partial autocorrelation of order  $h$  is

$$\pi(h) = \text{Corr}\left[X_t - \hat{X}_{h-1,t}^{(f)}, X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}\right].$$

In this section, we discuss partial autocorrelations in some detail. We show that  $\pi(h) = \phi_{h,h}$ , where  $\phi_{h,h}$  is the coefficient associated with  $X_{t-h}$  in the expression of the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ . As a by-product of the discussion, we also obtain an important recursion formula (known as *Durbin-Levinson* recursions) which relates the autoregressive coefficients of the  $AR(h)$  fit to those of  $AR(h-1)$  fit.

We first write down a basic result.

**Lemma 13.4.1.** *If  $\{X_t\}$  is a stationary  $AR(p)$  series, then  $\pi(h) = 0$  for  $h \geq p+1$ .*

It is fairly easy to see why this lemma is true. Assume that  $E[X_t] = 0$ . Since the series is assumed to be  $AR(p)$  with autoregressive coefficients  $\phi_1, \dots, \phi_p$ , for any  $h \geq p+1$ , the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h+1}$  is  $\hat{X}_{h-1,t}^{(f)} = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p}$  and  $X_t - \hat{X}_{h-1,t}^{(f)} = \varepsilon_t$ . If  $\hat{X}_{h-1,t-h}^{(b)}$  is the best linear predictor of  $X_{t-h}$  based on  $X_{t-h+1}, \dots, X_{t-1}$ , then the remainder  $X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}$  is a linear function of  $X_{t-1}, \dots, X_{t-h}$ . Since  $\varepsilon_t$  is uncorrelated with  $X_{t-1}, \dots, X_{t-h}$ , we can conclude that  $X_t - \hat{X}_{h-1,t} = \varepsilon_t$  is uncorrelated with  $X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}$  and thus  $\pi(h) = 0$ .

### 13.4.1 Expression for $\pi(h)$ and Durbin-Levinson Iterations

Let  $\{X_t\}$  be a mean zero stationary series which is not necessarily an autoregressive series of any finite order. Let  $\phi_{h,1}, \dots, \phi_{h,h}$  be the coefficients of the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$  (ie,  $\hat{X}_{h,t}^{(f)} = \phi_{h,1} X_{t-1} + \dots + \phi_{h,h} X_{t-h}$ ). The arguments here are given in terms of the theoretical autocovariances which can be replaced by their sample estimates for numerical computations based on the available data.

By Lemma 13.2.1

$$\begin{aligned}\hat{X}_{h-1,t}^{(f)} &= \phi_{h-1,1} X_{t-1} + \dots + \phi_{h-1,h-1} X_{t-h+1}, \\ \hat{X}_{h-1,t-h}^{(b)} &= \phi_{h-1,1} X_{t-h+1} + \dots + \phi_{h-1,h-1} X_{t-1},\end{aligned}$$

and the mean square error in forecasting  $X_t$  by  $\hat{X}_{h-1,t}^{(f)}$  is the same as the error in backcasting  $X_{t-h}$  by  $\hat{X}_{h-1,t-h}^{(b)}$ . This common prediction error is

$$\begin{aligned}PE(h-1) &= E\left[X_t - \hat{X}_{h,t}^{(f)}\right]^2 = E\left[X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}\right]^2 \\ &= \gamma(0) - [\phi_{h-1,1}\gamma(1) + \dots + \phi_{h-1,h-1}\gamma(h)].\end{aligned}$$

The identity  $\pi(h) = \phi_{h,h}$  and Durbin-Levinson iterative formula follow by equating the coefficients of two different but equivalent expressions of the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ . This basic idea is detailed in the following observations.

- (i) Any linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$  can be expressed as

$$L = a_1 X_{t-1} + \dots + a_{h-1} X_{t-h+1} + a_h \delta_{h-1,t-h},$$

where  $\delta_{h-1,t-h} = X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}$ .

- (ii) Since  $\delta_{h-1,t-h}$  is uncorrelated with  $X_{t-1}, \dots, X_{t-h+1}$ , the prediction error  $E[X_t - L]^2$ , where  $L$  is as in (i), is minimized at

$$\begin{aligned} a_1^* &= \phi_{h-1,1}, \dots, a_{h-1}^* = \phi_{h-1,h-1}, \text{ and} \\ a_h^* &= \text{Cov}[X_t, \delta_{h-1,t-h}] / \text{Var}[\delta_{h-1,t-h}] \\ &= \text{Cov}[X_t - \hat{X}_{h-1,t}^{(f)}, X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}] / \text{Var}[X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}] \\ &= \text{Corr}[X_t - \hat{X}_{h-1,t}^{(f)}, X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}] = \pi(h), \end{aligned}$$

where the last two steps follow from the fact that

$$E[X_t - \hat{X}_{h-1,t}^{(f)}]^2 = E[X_{t-h} - \hat{X}_{h-1,t-h}^{(b)}]^2.$$

- (iii) The coefficients of the best linear predictor  $\hat{X}_{h,t}^{(f)}$  are given in (ii) above, and using [Lemma 13.2.1](#) we have

$$\begin{aligned} \hat{X}_{h,t}^{(f)} &= a_1^* X_{t-1} + \dots + a_{h-1}^* X_{t-h+1} + a_h^* \delta_{h-1,t-h} \\ &= \phi_{h-1,1} X_{t-1} + \dots + \phi_{h-1,h-1} X_{t-h+1} \\ &\quad + \pi(h)[X_{t-h} - \phi_{h-1,1} X_{t-h+1} - \dots - \phi_{h-1,h-1} X_{t-1}] \\ &= [\phi_{h-1,1} - \pi(h)\phi_{h-1,h-1}] X_{t-1} + \dots \\ &\quad + [\phi_{h-1,h-1} - \pi(h)\phi_{h-1,1}] X_{t-h+1} + \pi(h) X_{t-h}. \end{aligned}$$

- (iv) The expression for the best linear predictor  $\hat{X}_{h,t}^{(f)}$  of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$  is

$$\phi_{h,1} X_{t-1} + \dots + \phi_{h,h-1} X_{t-h+1} + \phi_{h,h} X_{t-h},$$

and this should be the same as  $\hat{X}_{h,t}^{(f)}$  in (iii). Equating the coefficients of  $X_{t-1}, \dots, X_{t-h}$  in these two expressions, we have

$$\begin{aligned} \phi_{h,h} &= \pi(h) = a_h^* \\ &= \text{Cov}[X_t, \delta_{h-1,t-h}] / \text{Var}[\delta_{h-1,t-h}] \\ &= \frac{\gamma(h) - \phi_{h-1,1}\gamma(h-1) - \dots - \phi_{h-1,h-1}\gamma(1)}{\gamma(0) - \phi_{h-1,1}\gamma(1) - \dots - \phi_{h-1,h-1}\gamma(h-1)}, \text{ and} \\ \phi_{h,j} &= \phi_{h-1,j} - \phi_{h,h}\phi_{h-1,h-j}, \quad j = 1, \dots, h-1. \end{aligned}$$

From Step (iv), we have  $\phi_{h,h} = \pi(h)$ . The expressions in (iv), which relate the  $AR(h)$  coefficients to  $AR(h-1)$ , provide a recursion formula for calculating the autoregressive coefficients, which is known as the *Durbin-Levinson* algorithm. In order to obtain the estimate  $\hat{\pi}(h)$  of  $\pi(h)$ , it is enough to fit an  $AR(h)$  model and solve the Yule-Walker equations with  $\{\hat{\gamma}(j)\}$  in place of  $\{\gamma(j)\}$ . Also note that the Durbin-Levinson algorithm

provides an iterative scheme for solving the Yule-Walker equations starting from AR(1). We thus arrive at the following important results.

**Theorem 13.4.1.** *If  $\{X_t\}$  is a mean zero stationary series, and if  $\phi_{h,1}X_{t-1} + \cdots + \phi_{h,h}X_{t-h}$  is the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ , then the partial autocorrelation of order  $h \geq 2$  is given by  $\pi(h) = \phi_{h,h}$ . If  $\{X_t\}$  is stationary AR( $p$ ), then  $\sqrt{n}\hat{\pi}(h) \xrightarrow{\mathcal{D}} N(0, 1)$  as  $n \rightarrow \infty$  for  $h \geq p+1$ .*

**Theorem 13.4.2** (Durbin-Levinson Iterations). *For a mean zero stationary series  $\{X_t\}$ , if  $\phi_{h,1}X_{t-1} + \cdots + \phi_{h,h}X_{t-h}$  is the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h}$ , then the coefficients  $\phi_{h,1}, \dots, \phi_{h,h}$  are related to those of the best linear predictor of  $X_t$  based on  $X_{t-1}, \dots, X_{t-h-1}$  as given in observation (iv) above.*

We now provide a justification of the result on the asymptotic normality of  $\hat{\pi}(h)$ . We will see in [Section 13.8.1](#) that when  $\{X_t\}$  is AR( $p$ ), then  $\sqrt{n}(\hat{\phi}_h - \phi_h) \xrightarrow{\mathcal{D}} N_h(\mathbf{0}, \sigma^2 \boldsymbol{\Gamma}_h^{-1})$ ,  $h \geq p+1$ , where  $\phi_h = \boldsymbol{\Gamma}_h^{-1} \boldsymbol{\gamma}_h$  and  $\hat{\phi}_h$  is the solution of the Yule-Walker equations  $\hat{\boldsymbol{\Gamma}}_h \hat{\boldsymbol{\phi}} = \hat{\boldsymbol{\gamma}}_h$ . Thus,  $\sqrt{n}(\hat{\pi}(h) - \pi(h)) \xrightarrow{\mathcal{D}} N(0, \sigma^2 \boldsymbol{\Gamma}_h^{-1}(h, h))$ , where  $\boldsymbol{\Gamma}_h^{-1}(h, h)$  is the last element of the matrix  $\boldsymbol{\Gamma}_h^{-1}$ . From matrix algebra it is known that

- (i)  $\boldsymbol{\Gamma}_h^{-1}(h, h) = |\boldsymbol{\Gamma}_h| / |\boldsymbol{\Gamma}_{h-1}|$ , where  $|\cdot|$  denotes the determinant,
- (ii)  $\boldsymbol{\Gamma}_h = \begin{pmatrix} \boldsymbol{\Gamma}_{h-1} & \boldsymbol{\gamma}_{h-1} \\ \boldsymbol{\gamma}_{h-1}^T & \gamma(0) \end{pmatrix}$  and  $|\boldsymbol{\Gamma}_h| = |\boldsymbol{\Gamma}_{h-1}| [\gamma(0) - \boldsymbol{\gamma}_{h-1}^T \boldsymbol{\Gamma}_{h-1}^{-1} \boldsymbol{\gamma}_{h-1}]$ ,
- (iii)  $\boldsymbol{\Gamma}_h^{-1}(h, h) = 1 / [\gamma(0) - \boldsymbol{\gamma}_{h-1}^T \boldsymbol{\Gamma}_{h-1}^{-1} \boldsymbol{\gamma}_{h-1}]$ , and
- (iv) from [Lemma 13.2.1](#),  $\gamma(0) - \boldsymbol{\gamma}_{h-1}^T \boldsymbol{\Gamma}_{h-1}^{-1} \boldsymbol{\gamma}_{h-1} = E[X_t - X_{h-1,t}^{(f)}]^2$ .

If  $\{X_t\}$  is stationary AR( $p$ ), then  $X_{h-1,t}^{(f)} = X_{p,t}^{(f)}$  when  $h \geq p+1$  and

$$E[X_t - X_{h-1,t}^{(f)}]^2 = E[X_t - X_{p,t}^{(f)}]^2 = \gamma(0) - \boldsymbol{\gamma}_p^T \boldsymbol{\Gamma}_p^{-1} \boldsymbol{\gamma}_p,$$

Stationarity of  $\{X_t\}$  implies

$$\begin{aligned} \gamma(0) &= \text{Var}[X_t] = \text{Var}[\phi_1 X_{t-1} + \cdots + \phi_p X_{t-p} + \varepsilon_t] \\ &= \text{Var}[\phi_1 X_{t-1} + \cdots + \phi_p X_{t-p}] + \text{Var}[\varepsilon_t] \\ &= \boldsymbol{\phi}^T \boldsymbol{\Gamma}_p \boldsymbol{\phi} + \sigma^2 = \boldsymbol{\gamma}_p^T \boldsymbol{\Gamma}_p^{-1} \boldsymbol{\gamma}_p + \sigma^2, \text{ and hence} \\ \sigma^2 &= \gamma(0) - \boldsymbol{\gamma}_p^T \boldsymbol{\Gamma}_p^{-1} \boldsymbol{\gamma}_p, \text{ and} \end{aligned}$$

$$\sigma^2 \boldsymbol{\Gamma}_h^{-1}(h, h) = \sigma^2 / [\gamma(0) - \boldsymbol{\gamma}_{h-1}^T \boldsymbol{\Gamma}_{h-1}^{-1} \boldsymbol{\gamma}_{h-1}] = 1$$

for any  $h \geq p+1$ . This shows that  $\sqrt{n}\hat{\pi}(h) \xrightarrow{\mathcal{D}} N(0, 1)$  as  $n \rightarrow \infty$  for  $h \geq p+1$ .

## 13.5 Causality and Invertibility

This section addresses a number of important mathematical issues and theoretical results, but their proofs are not given. For theoretical details including proofs, the readers may consult the book by Brockwell and Davis [63].

A series  $\{X_t\}$  is called causal if it has an  $MA(\infty)$  representation, that is

$$X_t - \mu = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}, \text{ with } \psi_0 = 1, \quad (5)$$

where  $\{\varepsilon_t\}$  are mean zero iid with common variance  $\sigma^2$  and  $\sum |\psi_j| < \infty$ . Any causal series is stationary since for any  $h \geq 0$ ,

$$\text{Cov}[X_t, X_{t+h}] = \sigma^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+h} = \gamma(h),$$

depends only on  $h$ . In general, explicit expressions for these  $\psi$  weights are difficult to obtain (unless it is an  $MA(q)$  model) even though iterative formulas are available. Fortunately, packages such as R will calculate these  $\psi$  coefficients for an  $ARMA(p, q)$  model.

A series  $\{X_t\}$  is invertible if it can be written as an  $AR(\infty)$  series, that is, it has the representation

$$X_t - \mu = \sum_{j=1}^{\infty} \pi_j (X_{t-j} - \mu) + \varepsilon_t, \quad (6)$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and common variance  $\sigma^2$ . Clearly, any  $AR(p)$  model is invertible. However, an invertible model need not be stationary. For instance if  $X_t = X_{t-1} + \varepsilon_t$ , where  $\{\varepsilon_t\}$  are iid mean 0 with variance  $\sigma^2$ , then  $\{X_t\}$  is invertible but is not stationary ([Example 13.2.5](#)). Except for autoregressive models, there are no simple expressions for  $\{\pi_j\}$  for  $MA(q)$  or  $ARMA(p, q)$  models. However, a computing package such as R can be used to obtain them.

Mathematical conditions for invertibility and causality will be discussed later. An invertible expression is useful for obtaining the forecast values, whereas a causal representation makes it easy to obtain the variance of forecasts.

**Example 13.5.1.** A mean zero  $AR(1)$  series  $\{X_t\}$  may be rewritten as

$$\begin{aligned} X_t &= \phi X_{t-1} + \varepsilon_t = \phi(\phi X_{t-2} + \varepsilon_{t-1}) + \varepsilon_t \\ &= \phi^2 X_{t-2} + \phi \varepsilon_{t-1} + \varepsilon_t \\ &= \phi^2(\phi X_{t-3} + \varepsilon_{t-2}) + \phi \varepsilon_{t-1} + \varepsilon_t \\ &= \phi^3 X_{t-3} + \phi^2 \varepsilon_{t-2} + \phi \varepsilon_{t-1} + \varepsilon_t. \end{aligned}$$

We may repeat this argument to get

$$X_t = \phi^r X_{t-r} + \phi^{r-1} \varepsilon_{t-r+1} + \cdots + \phi \varepsilon_{t-1} + \varepsilon_t,$$

for any positive integer  $r$ . If  $|\phi| < 1$ , then  $\phi^r \rightarrow 0$  and hence  $\phi^r X_{t-r} \xrightarrow{P} 0$  as  $r \rightarrow \infty$ . Thus we may reexpress  $X_t$  as

$$X_t = \sum_{j=0}^{\infty} \phi^j \varepsilon_{t-j},$$

which is a causal representation of the series  $\{X_t\}$  with  $\psi_j = \phi^j$ .

**Example 13.5.2.** A mean zero MA(1) series may be rewritten as

$$\varepsilon_t = X_t - \theta \varepsilon_{t-1} = X_t - \theta(X_{t-1} - \theta \varepsilon_{t-2}) = X_t - \theta X_{t-1} + \theta^2 \varepsilon_{t-2}.$$

Repeating this argument we have, for any positive integer  $r$ ,

$$\varepsilon_t = X_t - \theta X_{t-1} + \cdots + (-\theta)^r X_{t-r} + (-\theta)^{r+1} \varepsilon_{t-r-1}.$$

If  $|\theta| < 1$ , then  $(-\theta)^{r+1} \varepsilon_{t-r-1} \xrightarrow{P} 0$  as  $r \rightarrow \infty$  and thus we have an invertible representation for the MA(1) series

$$\varepsilon_t = X_t - \theta X_{t-1} + \theta^2 X_{t-2} + \cdots, \text{ and hence}$$

$$X_t = \theta X_{t-1} - \theta^2 X_{t-2} + \cdots + \varepsilon_t = \sum_{j=1}^{\infty} \pi_j X_{t-j} + \varepsilon_t,$$

with  $\pi_j = (-1)^{j-1} \theta^j$ .

**Example 13.5.3.** Consider an ARMA(2, 2) model with  $\phi_1 = 0.8$ ,  $\phi_2 = -0.15$ ,  $\theta_1 = 0.6$ ,  $\theta_2 = 0.08$ . Using R, we have obtained the first 12 values of  $\{\psi_j\}$  starting with  $\psi_1 = 1.400$  are

$$[1.400, 1.050, 0.630, 0.347, 0.184, 0.094, 0.050, 0.024, 0.012, 0.006, 0.003, 0.002].$$

Note that the values of  $\psi_j$  become small for large  $j$  and this is typical for any ARMA series satisfying the condition of stationarity as will be seen later. Using R, we can get an invertible expression, and the first 12 values of  $\{\pi_j\}$  starting with  $\pi_1 = 1.400$  are

$$[1.400, -0.910, 0.434, -0.188, 0.078, -0.032, 0.013, -0.005, 0.002, -0.001, 0.0003, -0.0001].$$

The values of  $\pi_j$  rapidly approach 0 as  $j$  increases.

### 13.5.1 Usefulness of Invertible Representation

If we know the weights  $\{\pi_j\}$ , then it is easy to do the forecasting from the data  $\{X_1, \dots, X_n\}$  assuming that  $\pi_j \rightarrow 0$  as  $j \rightarrow \infty$ , which is generally true for ARMA processes under appropriate conditions. For notational simplicity let us assume that  $\mu = 0$ . Then we can write this series as

$$X_t = \varepsilon_t + \pi_1 X_{t-1} + \pi_2 X_{t-2} + \cdots.$$

If we have the entire past  $\{X_t: -\infty < t \leq n\}$ , then the forecasted value of  $X_{n+1}$  is

$$\hat{X}_{n+1} = \pi_1 X_n + \pi_2 X_{n-1} + \pi_3 X_{n-2} + \dots,$$

If  $X_{n+1}$  were known, then forecast of  $X_{n+2}$  would be

$$\pi_1 X_{n+1} + \pi_2 X_n + \pi_3 X_{n-1} + \dots.$$

Since  $X_{n+1}$  is unknown, then we can substitute it by  $\hat{X}_{n+1}$  (justification given in [Lemma 13.6.1](#)), leading to the forecasted value of  $X_{n+2}$  as

$$\hat{X}_{n+2} = \pi_1 \hat{X}_{n+1} + \pi_2 X_n + \pi_3 X_{n-1} + \dots.$$

This method can now be replicated to forecast  $X_{n+3}, X_{n+4}$ , etc.

Typically  $\mu$  is not equal to zero, but obtaining the forecasts is not difficult with the known  $\pi$  values. For instance, the forecasting formula for  $X_{n+1}$  is given by

$$\hat{X}_{n+1} - \mu = \pi_1(X_n - \mu) + \pi_2(X_{n-1} - \mu) + \pi_3(X_{n-2} - \mu) + \dots.$$

These expressions for forecasts assume that the entire past  $\{X_t: -\infty < t \leq n\}$  is known. However, if  $\pi_j \rightarrow 0$  rapidly as  $j \rightarrow \infty$ , which is the case for ARMA processes under the condition of invertibility, then the terms involving  $X_t$ ,  $t \leq 0$ , may be ignored since the associated  $\pi$ -coefficients are negligible and thus the approximate forecasts are linear functions of the available data  $X_1, \dots, X_n$ .

### 13.5.2 Usefulness of Causal Representation

Any practical approach to forecasting is incomplete without addressing the issue of prediction limits (or prediction intervals) of these forecasts. If the observed series is  $\{X_1, \dots, X_n\}$ , then the  $h$  step ahead forecast is denoted by  $\hat{X}_{n+h}$ . The forecast error is  $X_{n+h} - \hat{X}_{n+h}$ , which is not known since  $X_{n+h}$  is unknown. For all the cases we consider in this chapter, the mean of the forecast error is equal to zero (or close to zero when the parameters of the model are estimated). The mean square error of the forecast error is denoted by

$$\begin{aligned}\sigma^2(h) &= \text{Var}[(X_{n+h} - \hat{X}_{n+h})|X_1, \dots, X_n] \\ &= \text{E}[(X_{n+h} - \hat{X}_{n+h})^2|X_1, \dots, X_n].\end{aligned}$$

Thus a prediction interval for  $X_{n+h}$  with confidence coefficient  $1 - \alpha$  is  $\hat{X}_{n+1} \pm z_{\alpha/2} \sigma(h)$  if the series is assumed to be stationary Gaussian.

### 13.5.3 Important Technical issues

Here is a summary of the technical issues for  $AR(p)$ ,  $MA(q)$ , and  $ARMA(p, q)$  models which will be discussed in this section.

1. For  $AR(p)$  models, we need conditions on the autoregressive coefficients  $\phi_1, \dots, \phi_p$  in order to guarantee stationarity.

2. Moving average models are not in general identifiable. Thus for  $MA(q)$  models, we need conditions on the moving average coefficients  $\theta_1, \dots, \theta_q$  in order to guarantee identifiability.
3. For  $ARMA(p, q)$  models,
  - (i) the coefficients in the AR part must satisfy constraints to guarantee stationarity (as in (1) above),
  - (ii) the MA coefficients must satisfy constraints in order to guarantee identifiability of the model (as in (2) above),
  - (iii) conditions on the AR and the MA coefficients are needed in order to guarantee “nonredundancy.” This issue will be discussed in soon.

### Nonuniqueness of Moving Average Models

Moving average models are not unique. For instance, consider an  $MA(1)$  series  $X_t - \mu = \varepsilon_t + \theta\varepsilon_{t-1}$ , where  $\{\varepsilon_t\}$  is white noise with variance  $\sigma^2$ . Under the assumption of normality (ie,  $\{\varepsilon_t\}$  are normally distributed), any stationary series is completely characterized by the mean and autocovariances. Thus if two sequences have the same mean and autocovariance functions, they are equally good descriptions of the data, that is, they provide the same fit and they have the same predictive performances. Consider the following two models

$$X_t - \mu = \varepsilon_t + \theta\varepsilon_{t-1}, \quad X_t - \mu = \varepsilon'_t + (1/\theta)\varepsilon'_{t-1},$$

where  $\theta \neq 0$ ,  $\{\varepsilon_t\}$  is white noise with variance  $\sigma^2$ , and  $\{\varepsilon'_t\}$  is white noise with variance  $\theta^2\sigma^2$ . Note that we only observe the data  $\{X_t\}$ , not  $\{\varepsilon_t\}$  or  $\{\varepsilon'_t\}$ . Both models have the same mean  $\mu$ . All the autocovariances of lag 2 or higher are zero for both models. For the first model

$$\gamma(0) = (1 + \theta^2)\sigma^2, \quad \gamma(1) = \theta\sigma^2, \quad \text{and } 0 = \gamma(2) = \gamma(3) = \dots.$$

For the second model,

$$\begin{aligned} \gamma(0) &= [1 + (1/\theta)^2](\theta^2\sigma^2) = (1 + \theta^2)\sigma^2, \\ \gamma(1) &= (1/\theta)(\theta^2\sigma^2) = \theta\sigma^2, \quad \text{and } 0 = \gamma(2) = \gamma(3) = \dots. \end{aligned}$$

So both the models have identical mean and autocovariance structures. Hence they will provide identical fits and predictions. This nonuniqueness, or lack of identifiability, is problematic since there are multiple “correct” models.

How is this issue resolved? If the value of  $\theta$  is larger than 1 in magnitude and  $\text{Var}[\varepsilon_t] = \sigma^2$ , then we may as well consider the second model for which coefficient associated with  $\varepsilon_{t-1}$  is  $1/\theta$  whose magnitude is less than 1. Similarly if the moving average coefficient  $(1/\theta)$  of the second model is larger than 1 in magnitude, then we may decide to use the first model for which the coefficient would be smaller than 1 in magnitude. Thus we can always choose a model whose moving average coefficient is no larger than 1 in magnitude and this is what is done in practice.

The same issue of nonidentifiability comes up for general  $MA(q)$  models. One can restrict attention to those models with appropriate conditions on the MA parameters  $\theta_1, \dots, \theta_q$  needed for identifiability and this is what is usually done.

### Redundancy Issue for $ARMA(p, q)$ Models

If a series  $\{X_t\}$  is white noise (ie,  $X_t = \varepsilon_t$ , where  $\{\varepsilon_t\}$  is white noise), then subtracting  $0.5X_{t-1} = 0.5\varepsilon_{t-1}$  from this series we have

$$X_t - 0.5X_{t-1} = \varepsilon_t - 0.5\varepsilon_{t-1}, \quad \text{ie, } X_t = 0.5X_{t-1} + \varepsilon_t - 0.5\varepsilon_{t-1}.$$

Now it seems that the series  $\{X_t\}$  is  $ARMA(1, 1)$ , whereas in reality it is a white noise. As a matter of fact we can rewrite  $X_t$  as

$$X_t = \phi X_{t-1} + \varepsilon_t - \phi \varepsilon_{t-1},$$

for any  $-1 < \phi < 1$ . Once again it looks as if  $\{X_t\}$  is  $ARMA(1, 1)$ ; however, there is a redundant parameter  $\phi$ . Also note that the number of such redundant models is infinite. In general if  $\theta = -\phi$  in an  $ARMA(1, 1)$  model  $X_t = \phi X_{t-1} + \varepsilon_t + \theta \varepsilon_{t-1}$ , then there is redundancy. When  $\theta \neq \phi$  then this redundancy is no longer present. The same issue needs to be addressed for general ARMA models and constraints on the parameters are needed in order to avoid any redundancy.

### Condition for Stationarity for $AR(p)$ Models

If  $\{X_t\}$  follows an  $AR(p)$  model, then it is stationary (and causal) if it can be written in the form Eq. (5). Conditions on the autoregressive coefficients  $\phi_1, \dots, \phi_p$  are needed to guarantee that the series can be written in this form. Let  $\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p$  be a polynomial in  $z$ , where  $z$  is complex. This polynomial has  $p$  roots which can be real or complex valued.

**Lemma 13.5.1.** *If the absolute values of all the roots of the polynomial  $\phi(z)$  of an  $AR(p)$  series are larger than 1, then the series is stationary.*

When  $p = 1$ , the root of the polynomial is  $1/\phi_1$ . The condition that the absolute value of  $1/\phi_1$  is larger than 1 is equivalent to the condition  $-1 < \phi_1 < 1$ . In this case, the autocorrelation function is  $\rho(h) = \phi_1^h$ ,  $h = 0, 1, \dots$ , and  $\rho(h)$  converges to zero exponentially fast as  $h$  increases (Example 13.2.3).

For  $p = 2$ , the condition on  $\phi_1$  and  $\phi_2$  for stationarity is a bit more complicated. The condition is:  $-1 < \phi_2 < 1$  and  $-1 < \phi_1/(1 - \phi_2) < 1$ , that is,  $(\phi_1, \phi_2)$  is inside the triangle  $\Delta = \{(u_1, u_2): -1 < u_2 < 1, -1 < u_1/(1 - u_2) < 1\}$ . The roots of  $\phi(z)$  are reciprocals of the roots of  $g(z) = z^2 - \phi_1 z - \phi_2$ . The roots of  $g(z)$  are  $(1/2) \left[ \phi_1 \pm \sqrt{\phi_1^2 + 4\phi_2} \right]$ , and they must be smaller than 1 in magnitude for the  $AR(2)$  series to be stationary. In the  $AR(2)$  case, it can be shown that  $\rho(h) = \phi_1 \rho(h-1) + \phi_2 \rho(h-2)$ ,  $h = 2, 3, \dots$ . When the roots  $r_1$  and  $r_2$  of  $g(z)$  are real and distinct, the theory of difference equations tells us that the autocorrelations  $\rho(h)$  behave like  $c_1 r_1^{-h} + c_2 r_2^{-h}$ , where  $c_1$  and  $c_2$  are real. When  $r_1 = r_2$ , that

is,  $\phi_1^2 + 4\phi_2 = 0$ , then  $\rho(h) = c_1 r_1^h + c_2 h r_1^h$ ,  $c_1, c_2$  real. Thus the autocorrelation function  $\{\rho(h)\}$  decays exponentially in  $h$  when the roots of  $g(z)$  are real.

When  $r_1$  and  $r_2$  are complex, the AR(2) can model series with pseudo-cyclical behavior. In this case,  $\phi_1^2 + 4\phi_2 < 0$ ,  $r_2$  is the complex conjugate of  $r_1$  and  $\rho(h)$  is of the form  $c_1 r_1^h + c_2 r_2^h$ , where  $c_1$  and  $c_2$  are complex. Calculations show that  $\rho(h) = |r_1|^h [\text{sign}(\phi_1)]^h \sin(2\pi Ch + D)/\sin(D)$ , where  $\text{sign}(\phi_1) = 1$  if  $\phi_1 > 0$  and  $= -1$  if  $\phi_1 < 0$ , and  $C$  and  $D$  are constants. Since  $|r_1| = |r_2| < 1$ ,  $\rho(h)$  decays exponentially in a sinusoidal fashion.

In general, the autocorrelations of an AR( $p$ ) process decay as a mixture of exponentials or as damped (sinusoidal) exponentials.

### Identifiability of MA( $q$ ) Models

Consider the polynomial  $\theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q$ , where  $\theta_1, \dots, \theta_q$  are the parameters of the MA( $q$ ) model. This polynomial has  $q$  roots which can be real or complex valued.

**Lemma 13.5.2.** *If the absolute values of all the roots of the polynomial  $\theta(z)$  of an MA( $q$ ) series are larger than or equal to 1, then the series is identifiable.*

For the MA(1) model, this translates into the condition that  $-1 \leq \theta_1 \leq 1$ . For MA(2) this condition is equivalent to the condition that  $(-\theta_1, -\theta_2)$  is in the triangle  $\Delta$  as in the AR(2) case, except that  $(-\theta_1, -\theta_2)$  is allowed to be on the boundaries of the triangle.

**Lemma 13.5.3.** *If for an MA( $q$ ) or ARMA( $p, q$ ) series, the absolute values of all the roots of the polynomial  $\theta(z)$  are larger than 1, then the series is invertible.*

For the MA(1) model, the conditions for identifiability and invertibility are  $-1 \leq \theta_1 \leq 1$  and  $-1 < \theta_1 < 1$ , respectively. For MA(2) the model is invertible if  $(-\theta_1, -\theta_2)$  is in the triangle  $\Delta$  as in the AR(2) case. Thus the condition for invertibility guarantees identifiability.

### Stationarity, Invertibility, and Nonredundancy of ARMA( $p, q$ ) Models

We have already seen that for an ARMA(1, 1) series, there is no redundancy when  $\theta \neq -\phi$ . For this case the roots of  $\phi(z)$  and  $\theta(z)$  are  $1/\phi$  and  $-1/\theta$ , respectively. Thus nonredundancy is achieved when  $\phi(z)$  and  $\theta(z)$  have no common root. This same condition is true for the general ARMA model.

**Lemma 13.5.4.** *For an ARMA( $p, q$ ) series, if  $\phi(z)$  and  $\theta(z)$  have no common root, then the model is nonredundant.*

For the general ARMA model we need the following:

- (a) roots of  $\phi(z)$  are larger than 1 in magnitude (condition for stationarity),
- (b) roots of  $\theta(z)$  are larger than or equal to 1 (condition for identifiability), and
- (c) the roots of  $\phi(z)$  are distinct from the roots of  $\theta(z)$  (condition for nonredundancy).

These are summarized in the following result.

**Lemma 13.5.5.** *An ARMA( $p, q$ ) series is stationary, identifiable, and nonredundant if the following conditions hold:*

- (a) the roots of  $\phi(z)$  are larger than 1 in magnitude,
- (b) the roots of  $\theta(z)$  are larger than or equal to 1 in magnitude,
- (c)  $\phi(z)$  and  $\theta(z)$  have no common roots.

For invertibility, in addition to stationarity and nonredundancy, condition (b) needs to be replaced by the stronger condition

- (b') the roots of  $\theta(z)$  are larger than 1 in magnitude.

## 13.6 Forecasting

This section deals with the issues of forecasting, and obtaining prediction intervals for an ARMA series  $\{X_t\}$ . Here the parameters associated with the ARMA model are assumed to be known, but they need to be estimated in practice, and the issue of estimation will be discussed later. Suppose that observations are  $X_1, \dots, X_n$  and we wish to forecast  $X_{n+1}, X_{n+2}, \dots$ , then a linear predictor of  $X_{n+h}$  is of the form  $a_0 + a_1 X_1 + \dots + a_n X_n$ . If all the variables are centered (ie, they are subtracted by their means), then we may take  $a_0 = 0$ . For notational convenience, we consider the issue of forecasting a centered series and we continue to denote the centered series as  $\{X_t\}$ .

If we wish to predict  $X_{n+1}$  using a linear predictor based on  $\mathbf{X}_n = (X_1, \dots, X_n)^T$  for the AR(1) case which is modeled as  $X_{n+1} = \phi X_n + \varepsilon_{n+1}$ , the best linear predictor of  $X_{n+1}$  is  $\hat{X}_{n+1} = \phi X_n$ . If  $X_{n+1}$  were available, then the best linear predictor of  $X_{n+2}$  based on  $\mathbf{X}_{n+1} = (X_1, \dots, X_n, X_{n+1})$  would be  $\phi X_{n+1}$ . Now  $X_{n+1}$  is unobserved, but a predicted value  $\hat{X}_{n+1}$  of  $X_{n+1}$  is available. Thus substituting  $\hat{X}_{n+1}$  in place of  $X_{n+1}$ , we can obtain a linear predictor of  $X_{n+2}$  based on  $\mathbf{X}_n$  as  $\hat{X}_{n+2} = \phi \hat{X}_{n+1} = \phi^2 X_n$ . We can similarly obtain a linear predictor of  $X_{n+h}$  based on  $X_1, \dots, X_n$  as  $\hat{X}_{n+h} = \phi^h X_n$ . Is  $\hat{X}_{n+h} = \phi^h X_n$  the best linear predictor of  $X_{n+h}$  based on  $X_1, \dots, X_n$ ? The answer is yes.

Let  $\hat{X}_{n+h}$  denote the best linear predictor of  $X_{n+h}$  based on  $X_1, \dots, X_n$ . In general, it is easy to obtain an expression for the best predictor of  $X_{n+h}$  based on  $X_1, \dots, X_{n+h-1}$  and in this expression we can substitute  $\hat{X}_{n+1}, \dots, \hat{X}_{n+h-1}$  for  $X_{n+1}, \dots, X_{n+h-1}$  in order to obtain the best linear predictor of  $X_{n+h}$  based on  $X_1, \dots, X_n$ . This is based on a rather simple argument as outlined below.

If  $Y$  is a rv and  $\mathbf{W}_1, \mathbf{W}_2$  are two random vectors, then the best predictors of  $Y$  given  $\mathbf{W}_1$ , and given  $\mathbf{W}_2$  and  $\mathbf{W}_1$  are  $E[Y|\mathbf{W}_1]$  and  $E[Y|\mathbf{W}_2, \mathbf{W}_1]$ , respectively. The law of iterative expectations tells us  $E[Y|\mathbf{W}_1] = E[E[Y|\mathbf{W}_2, \mathbf{W}_1]|\mathbf{W}_1]$ . Is this result true for best linear predictors? The answer is yes if  $Y, \mathbf{W}_1, \mathbf{W}_2$  are jointly normally distributed, since the best predictor is the same as the best linear predictor. However, the result is true more generally in the sense described in the result below. Let the best linear predictors of  $Y$  given  $\mathbf{W}_1$ , and  $Y$  given  $\mathbf{W}_1, \mathbf{W}_2$  be denoted by  $L(Y|\mathbf{W}_1)$  and  $L(Y|\mathbf{W}_2, \mathbf{W}_1)$ , respectively. We can then write the following result.

**Lemma 13.6.1.** *Let  $L(\cdot|\cdot)$  be the best linear predictor as described above.*

- (a)  $L$  is linear in the response in the sense that for any rv's  $Y_1$  and  $Y_2$ ,

$$L(Y_1 + Y_2 | \mathbf{W}_1) = L(Y_1 | \mathbf{W}_1) + L(Y_2 | \mathbf{W}_1).$$

- (b)  $L$  satisfies the iterative formula  $L(Y|\mathbf{W}_1) = L[L(Y|\mathbf{W}_2, \mathbf{W}_1)|\mathbf{W}_1]$ .  
(c) If  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are uncorrelated, that is,  $\text{Cov}[\mathbf{W}_1, \mathbf{W}_2] = \mathbf{0}$ , then

$$L[L(Y|\mathbf{W}_2)|\mathbf{W}_1] = 0 \quad \text{and} \quad L(Y|\mathbf{W}_2, \mathbf{W}_1) = L(Y|\mathbf{W}_2) + L(Y|\mathbf{W}_1).$$

This proof of this lemma is left as an exercise. It can be used to justify replacing  $X_{n+1}, \dots, X_{n+h-1}$  by  $\hat{X}_{n+1}, \dots, \hat{X}_{n+h-1}$  in the expression for the best linear predictor of  $X_{n+h}$  given  $X_1, \dots, X_{n+h-1}$  in order to obtain the best linear predictor of  $X_{n+h}$  given  $X_1, \dots, X_n$ .

### 13.6.1 Forecasting an $AR(p)$ Series

Forecasting with an  $AR(p)$  model with autoregressive coefficients  $\phi_1, \dots, \phi_p$  is quite simple as it has a regression form. Suppose that observations are  $X_1, \dots, X_n$  and we wish to forecast  $X_{n+1}, X_{n+2}, \dots$ . The best linear predictor of  $X_{n+h}$  based on  $X_1, \dots, X_n$  will be denoted by  $\hat{X}_{n+h}$ . Since

$$X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \varepsilon_t,$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ , the best linear predictor of  $X_{n+1}$  based on  $X_1, \dots, X_n$  is

$$\hat{X}_{n+1} = \phi_1 X_n + \dots + \phi_p X_{n+1-p}.$$

If  $X_{n+1}$  were known the best linear predictor based on  $X_1, \dots, X_{n+1}$  is  $\phi_1 X_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p}$ . In [Lemma 13.6.1](#), use  $Y = X_{n+2}$ ,  $\mathbf{W}_1 = (X_1, \dots, X_n)^T$  and  $\mathbf{W}_2 = X_{n+1}$  to obtain

$$\begin{aligned}\hat{X}_{n+2} &= L(Y|\mathbf{W}_1) = L[L(Y|\mathbf{W}_2, \mathbf{W}_1)|\mathbf{W}_1] \\ &= L[\phi_1 X_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p}|\mathbf{W}_1] \\ &= \phi_1 L(X_{n+1}|\mathbf{W}_1) + \phi_2 X_n + \dots + \phi_p X_{n+2-p} \\ &= \phi_1 \hat{X}_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p}.\end{aligned}$$

We have used linearity of  $L$  (part (a) of [Lemma 13.6.1](#)) and the fact that  $L(X_t|\mathbf{W}_1) = X_t$  for any  $t = 1, \dots, n$ .

If we denote  $\hat{X}_t = X_t$ ,  $t = 1, \dots, n$ , then the argument used above can be employed to show that for any  $h \geq 1$ ,

$$\hat{X}_{n+h} = \phi_1 \hat{X}_{n+h-1} + \dots + \phi_p \hat{X}_{n+h-p}.$$

### 13.6.2 Forecasting an $MA(q)$ Series

For a mean zero  $MA(q)$  series with coefficients  $\theta_1, \dots, \theta_q$ ,  $X_{n+h}$  is uncorrelated with  $X_n$  when  $h \geq q+1$ . Thus  $\hat{X}_{n+h} = 0$  for  $h = q+1, q+2, \dots$ . We discuss below how to find the formula for forecasting  $\hat{X}_{n+1}, \dots, \hat{X}_{n+q}$ . As will be clear in the subsequent discussion, unlike in the  $AR(p)$  case where the forecasted value  $\hat{X}_{n+h}$  depends only on  $X_n, \dots, X_{n+1-p}$ , here the forecasted value of  $X_{n+h}$ ,  $1 \leq h \leq q$ , depends on the entire available past  $X_n, \dots, X_1$ .

For an invertible series given in Eq. (6), the forecasting formula is simple if the entire past  $\{X_t: -\infty < t \leq n\}$  were available. Absolute summability of  $\{\pi_j\}$  guarantees that  $\pi_j \rightarrow 0$  as  $j \rightarrow \infty$ , and for an invertible  $MA(q)$  model,  $\pi_j$  decays exponentially fast as  $j$  increases. So when predicting  $X_{n+h}$  based on  $X_1, \dots, X_n$ , we may simply truncate the  $\pi$  series at  $j = n$  thus approximating the process by an  $AR(p)$  sequence with  $p = n$ , and then use the methods associated with forecasting an autoregressive process as given in [Section 13.6.1](#).

We now discuss how to carry out forecasting without having to obtain the values of  $\{\pi_j\}$  assuming that the series is invertible. Note that

$$X_{n+1} = \varepsilon_{n+1} + \theta_1 \varepsilon_n + \dots + \theta_q \varepsilon_{n+1-q}.$$

For the moment assume that in addition to the observations  $\mathbf{X}_n = (X_1, \dots, X_n)^T$  we also have  $\boldsymbol{\varepsilon}_0 = (\varepsilon_{-q+1}, \dots, \varepsilon_0)^T$ , and we predict  $X_{n+1}$  based on  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X}_n$ . Any linear combination of  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X}_n$  can be rewritten as a linear combination of  $\varepsilon_{-q+1}, \dots, \varepsilon_n$ , and vice versa. Thus the best linear predictor of  $X_{n+1}$  based on  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X}_n$  is given by

$$\hat{X}_{n+1} = \theta_1 \varepsilon_n + \dots + \theta_q \varepsilon_{n+1-q}.$$

In order to forecast  $X_{n+2}$ , note that

$$X_{n+2} = \varepsilon_{n+2} + \theta_1 \varepsilon_{n+1} + \theta_2 \varepsilon_n + \dots + \theta_q \theta_{n+2-q},$$

and  $(\varepsilon_{n+2}, \varepsilon_{n+1})$  is uncorrelated with  $\mathbf{W}_1 = (\boldsymbol{\varepsilon}_0^T, X_1, \dots, X_n)^T$ . Thus the best linear predictor of  $X_{n+2}$  given  $\mathbf{W}_1$  is

$$\hat{X}_{n+2} = \theta_2 \varepsilon_n + \dots + \theta_q \theta_{n+2-q}.$$

A similar argument will show that

$$\begin{aligned}\hat{X}_{n+3} &= \theta_3 \varepsilon_n + \dots + \theta_q \varepsilon_{n+3-q}, \\ \hat{X}_{n+h} &= \theta_h \varepsilon_n + \dots + \theta_q \varepsilon_{n+h-q}, \quad 1 \leq h \leq q, \text{ and} \\ \hat{X}_{n+h} &= 0, \quad h > q.\end{aligned}$$

The forecasts  $\hat{X}_{n+h}$ ,  $h = 1, \dots, q$ , depend on knowing  $\varepsilon_n, \varepsilon_{n-1}, \dots, \varepsilon_{n+1-q}$ . We now point out how to obtain these from  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X}_n$ . We may obtain  $\varepsilon_1$  and  $\varepsilon_2$  as

$$\begin{aligned}\varepsilon_1 &= X_1 - (\theta_1 \varepsilon_0 + \dots + \theta_q \varepsilon_{1-q}), \text{ and} \\ \varepsilon_2 &= X_2 - (\theta_1 \varepsilon_1 + \dots + \theta_q \varepsilon_{2-q}).\end{aligned}$$

Continuing this way, once we have  $\varepsilon_1, \dots, \varepsilon_t$ , we may obtain

$$\varepsilon_{t+1} = X_{t+1} - (\theta_1 \varepsilon_t + \dots + \theta_q \varepsilon_{t+1-q}), \quad t = 1, \dots, n-1.$$

It is important to note that only the data  $\mathbf{X}_n = (X_1, \dots, X_n)$  is available and not  $\boldsymbol{\varepsilon}_0$ . Even though  $\varepsilon_n, \dots, \varepsilon_{n+1-q}$  are linear combinations of  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X}_n$ , often one takes  $\boldsymbol{\varepsilon}_0 = \mathbf{0}$  since the coefficients associated with  $\boldsymbol{\varepsilon}_0$  in these linear combinations are negligible if the assumption of invertibility is valid. Without the assumption of invertibility, it is still possible to obtain the forecasts, but that issue will not be discussed here.

**Example 13.6.1.** Consider a mean zero MA(1) series with the moving average coefficient  $\theta$ . Assuming  $\varepsilon_0$  to be known, we can get

$$\varepsilon_1 = X_1 - \theta\varepsilon_0, \quad \varepsilon_2 = X_2 - \theta\varepsilon_1, \dots, \quad \varepsilon_n = X_n - \theta\varepsilon_{n-1}.$$

Using the arguments given in [Example 13.5.2](#), we have

$$\varepsilon_n = X_n - \theta X_{n-1} + \theta^2 X_{n-2} + \dots + (-1)^{n-1} \theta^{n-1} X_1 + (-1)^n \theta^n \varepsilon_0.$$

The forecasted value  $\hat{X}_{n+1}$  of  $X_{n+1}$  based on  $\varepsilon_0, X_1, \dots, X_n$  is

$$\begin{aligned}\hat{X}_{n+1} &= \theta\varepsilon_n \\ &= \theta[X_n - \theta X_{n-1} + \theta^2 X_{n-2} + \dots + (-1)^{n-1} \theta^{n-1} X_1] + (-1)^n \theta^{n+1} \varepsilon_0.\end{aligned}$$

The term  $\theta^{n+1} \varepsilon_0$  is negligible if  $|\theta| < 1$  (ie, the series is invertible) and thus  $\hat{X}_{n+1}$  is approximately a linear combination of  $X_1, \dots, X_n$ .

### 13.6.3 Forecasting an ARMA Series

If  $\{X_t\}$  is mean zero  $ARMA(p, q)$ , then the method for forecasting  $X_{n+h}$  based on  $\mathbf{X}_n = (X_1, \dots, X_n)^T$  combines the methods given for AR and MA models. If the series is invertible, then we can obtain, in principle, the (approximate) best linear predictors of  $X_{n+h}$ ,  $h = 1, 2, \dots$ , based on  $X_1, \dots, X_n$ , by approximating  $\{X_t\}$  by an  $AR(p)$  process with  $p = n$ .

For the moment assume that  $\boldsymbol{\varepsilon}_p = (\varepsilon_{p+1-q}, \dots, \varepsilon_p)^T$  is known. Then we can obtain  $\varepsilon_{p+1}, \dots, \varepsilon_n$  as linear combinations of  $\boldsymbol{\varepsilon}_p$  and  $\mathbf{X}_n$  as will be shown below.

Since

$$\begin{aligned}X_{n+1} &= \phi_1 X_n + \dots + \phi_p X_{n+1-p} + \varepsilon_{n+1} + \theta_1 \varepsilon_n + \dots + \theta_q \varepsilon_{n+1-q} \\ &= \phi_1 X_n + \dots + \phi_p X_{n+1-p} + \theta_1 \varepsilon_n + \dots + \theta_q \varepsilon_{n+1-q} + \varepsilon_{n+1},\end{aligned}$$

the best linear predictor  $\hat{X}_{n+1}$  of  $X_{n+1}$  is given by

$$\hat{X}_{n+1} = \phi_1 X_n + \dots + \phi_p X_{n+1-p} + \theta_1 \varepsilon_n + \dots + \theta_q \varepsilon_{n+1-q}.$$

Since

$$\begin{aligned}X_{n+2} &= \phi_1 X_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p} \\ &\quad + \theta_1 \varepsilon_{n+1} + \theta_2 \varepsilon_n + \dots + \theta_q \varepsilon_{n+2-q} + \varepsilon_{n+2},\end{aligned}$$

the best linear predictor of  $X_{n+2}$  would be

$$\phi_1 X_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p} + \theta_1 \varepsilon_{n+1} + \theta_2 \varepsilon_n + \dots + \theta_q \varepsilon_{n+2-q},$$

if  $X_{n+1}, \varepsilon_{n+1}$  were known. Following [Lemma 13.6.1](#) we can now replace  $X_{n+1}$  by  $\hat{X}_{n+1}$  and  $\varepsilon_{n+1}$  by  $\hat{\varepsilon}_{n+1} = 0$  (since  $\varepsilon_{n+1}$  is uncorrelated with  $\boldsymbol{\varepsilon}_p$  and  $X_1, \dots, X_n$ ). The best linear predictor of  $X_{n+2}$  given  $\boldsymbol{\varepsilon}_p$  and  $X_1, \dots, X_n$  is

$$\hat{X}_{n+2} = \phi_1 \hat{X}_{n+1} + \phi_2 X_n + \dots + \phi_p X_{n+2-p} + \theta_2 \varepsilon_n + \dots + \theta_q \varepsilon_{n+2-q}.$$

A similar argument will show that

$$\begin{aligned}\hat{X}_{n+h} &= \phi_1 \hat{X}_{n+h-1} + \cdots + \phi_p \hat{X}_{n+h-p} + \theta_h \varepsilon_n + \cdots + \theta_q \varepsilon_{n+h-q}, \quad h = 1, \dots, q, \\ \hat{X}_{n+h} &= \phi_1 \hat{X}_{n+h-1} + \cdots + \phi_p \hat{X}_{n+h-p}, \quad h > q,\end{aligned}$$

with the understanding that  $\hat{X}_t = X_t$  for  $t = 1, \dots, n$ .

The forecasts  $\hat{X}_{n+h}$ ,  $h \geq 1$ , depend on  $\varepsilon_{n+1-q}, \dots, \varepsilon_n$ . We need to obtain their values when  $\varepsilon_p$  and  $X_1, \dots, X_n$  are available, and this can be done iteratively starting with  $\varepsilon_{p+1}$

$$\begin{aligned}\varepsilon_{p+1} &= X_{p+1} - \phi_1 X_p - \cdots - \phi_p X_1 - (\theta_1 \varepsilon_p + \cdots + \theta_q \varepsilon_{p+1-q}), \\ \varepsilon_{p+2} &= X_{p+2} - \phi_1 X_{p+1} - \cdots - \phi_p X_2 - (\theta_1 \varepsilon_{p+1} + \cdots + \theta_q \varepsilon_{p+2-q}),\end{aligned}$$

and, when  $\varepsilon_{p+1}, \dots, \varepsilon_t$ ,  $t \geq p+1$ , are available, we can find

$$\varepsilon_{t+1} = X_{t+1} - \phi_1 X_t - \cdots - \phi_p X_{t+1-p} - (\theta_1 \varepsilon_t + \cdots + \theta_q \varepsilon_{t+1-q}).$$

Thus  $\varepsilon_n, \dots, \varepsilon_{n+1-q}$  and hence  $\hat{X}_{n+h}$ ,  $h \geq 1$ , are linear combinations of  $\varepsilon_p$  and  $X_1, \dots, X_n$ . Under the conditions of invertibility, the coefficients associated with  $\varepsilon_p, \dots, \varepsilon_{p+1-q}$  are negligible. In practice, often  $\varepsilon_p, \dots, \varepsilon_{p+1-q}$  are taken to be zeros.

**Example 13.6.2.** Let us assume that we have an ARMA(1, 1) process  $X_t = \phi X_{t-1} + \varepsilon_t + \theta \varepsilon_{t-1}$ ,  $\phi \neq -\theta$ , and we want to predict  $X_{n+1}, X_{n+2}, \dots$  using the observations  $X_1, \dots, X_n$ . If  $\varepsilon_1$  were available, then generate  $\varepsilon_2, \varepsilon_3, \dots, \varepsilon_n$  as

$$\varepsilon_2 = X_2 - (\phi X_1 + \theta \varepsilon_1), \quad \varepsilon_3 = X_3 - (\phi X_2 + \theta \varepsilon_2), \dots, \quad \varepsilon_n = X_n - (\phi X_{n-1} + \theta \varepsilon_{n-1})$$

Thus predicted values of  $X_{n+h}$ ,  $h \geq 1$ , are

$$\begin{aligned}\hat{X}_{n+1} &= \phi X_n + \theta \varepsilon_n, \quad \hat{X}_{n+2} = \phi \hat{X}_{n+1} + \theta \hat{\varepsilon}_{n+1} = \phi \hat{X}_{n+1}, \text{ and} \\ \hat{X}_{n+h+1} &= \phi \hat{X}_{n+h}, \quad h \geq 2.\end{aligned}$$

*Remark 13.6.1.* For the MA( $q$ ) and ARMA( $p, q$ ) cases, it is generally assumed that the series is invertible. However, the condition of invertibility is more of a convenience than necessity. Even without the condition of invertibility, it is still possible to forecast  $X_{n+h}$ ,  $h \geq 1$ , using the best linear predictor  $a_1 X_n + \cdots + a_{n-1} X_1$ , but unlike in the invertible case, the coefficients  $\{a_j\}$  may not converge to 0 rapidly. For instance, for the mean zero MA(1) case if we have  $\theta = -1$ , then the best linear predictor of  $X_{n+1}$  is of the form

$$\hat{X}_{n+1} = \sum_{j=0}^{n-1} a_j X_{n-j}, \text{ where } a_j = -1 + (j+1)/(n+1).$$

Note that the coefficients  $\{a_j\}$  increase linearly from  $a_0 = -n/(n+1)$  to  $a_{n-1} = -1/(n+1)$ .

### 13.6.4 Standard Error for Prediction

Calculations of standard errors of predictions are needed to construct confidence bounds. Let  $\sigma^2(h) = E[(X_{n+h} - \hat{X}_{n+h})^2 | \mathbf{X}_n]$  be the mean square error for predicting  $X_{n+h}$  based on the data  $\mathbf{X}_n = (X_1, \dots, X_n)^T$ . Once again we assume that the parameters of the ARMA model are known and formulas for  $\sigma^2(h)$  depend on the unknown parameters. In practice these parameters are replaced by their sample estimates. Under the assumption of normality, a prediction interval for  $X_{n+h}$  with confidence coefficient  $1 - \alpha$  is given by  $\hat{X}_{n+h} \pm z_{\alpha/2}\sigma(h)$ .

The formulas for  $\sigma^2(h)$  can be obtained easily for a series with causal representation

$$X_t = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}, \text{ with } \psi_0 = 1,$$

where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ , and  $\sum |\psi_j| < \infty$ . If we assume that the entire past  $\mathbf{X}_n = \{X_t: -\infty < t \leq n\}$  is known, then it is equivalent to knowing the entire past of  $\{\varepsilon_t: -\infty < t \leq n\}$  of the errors. If a predictor of  $X_{n+h}$  is a linear combination of  $\{X_t: -\infty < t \leq n\}$ , then it is a linear combination of  $\{\varepsilon_t: -\infty < t \leq n\}$ . Under appropriate conditions of stationarity and invertibility, the dependence on  $\{\varepsilon_t: -\infty < t \leq 0\}$  is negligible. Since

$$X_{n+1} = \sum_{j=1}^{\infty} \psi_j \varepsilon_{n+1-j} + \varepsilon_{n+1},$$

the best linear predictor  $\hat{X}_{n+1}$  of  $X_{n+1}$  is  $\sum_{j=1}^{\infty} \psi_j \varepsilon_{n+1-j}$ , and the mean square error of prediction is

$$\sigma^2(1) = E\left[\left(X_{n+1} - \hat{X}_{n+1}\right)^2 | \mathbf{X}_n\right] = E\left[\varepsilon_{n+1}^2\right] = \sigma^2.$$

The best linear predictor of  $X_{n+2}$  given  $\mathbf{X}_n$  is

$$\begin{aligned} \hat{X}_{n+2} &= \psi_2 \varepsilon_n + \psi_3 \varepsilon_{n-1} + \dots, \text{ and} \\ X_{n+2} - \hat{X}_{n+2} &= [\varepsilon_{n+2} + \psi_1 \varepsilon_{n+1} + \psi_2 \varepsilon_n + \dots] \\ &\quad - [\psi_2 \varepsilon_n + \psi_3 \varepsilon_{n-1} + \dots] \\ &= \varepsilon_{n+2} + \psi_1 \varepsilon_{n+1}. \end{aligned}$$

Hence the mean square error for predicting  $X_{n+2}$  by  $\hat{X}_{n+2}$  is

$$\begin{aligned} \sigma^2(2) &= E\left[\left(X_{n+2} - \hat{X}_{n+2}\right)^2 | \mathbf{X}_n\right] = E\left[\left(\varepsilon_{n+2} + \psi_1 \varepsilon_{n+1}\right)^2 | \mathbf{X}_n\right] \\ &= \sigma^2 + \psi_1^2 \sigma^2 = (1 + \psi_1^2) \sigma^2. \end{aligned}$$

A very similar argument will show that the best linear predictor of  $X_{n+h}$  based on the entire past and the mean square error of prediction are

$$\begin{aligned} \hat{X}_{n+h} &= \psi_h \varepsilon_n + \psi_{h+1} \varepsilon_{n-1} + \dots, \\ X_{n+h} - \hat{X}_{n+h} &= \varepsilon_{n+h} + \psi_1 \varepsilon_{n+h-1} + \dots + \psi_{h-1} \varepsilon_{n+1}, \text{ and} \end{aligned}$$

$$\begin{aligned}\sigma^2(h) &= \text{E}\left[\left(X_{n+h} - \hat{X}_{n+h}\right)^2 | \mathbf{X}_n\right] \\ &= (1 + \psi_1^2 + \cdots + \psi_{h-1}^2)\sigma^2.\end{aligned}$$

Note that  $\sigma^2(h) \rightarrow \sigma^2 \sum_{j=0}^{\infty} \psi_j^2$  as  $h \rightarrow \infty$ , where  $\psi_0 = 1$ . The limit of  $\sigma^2(h)$  is a constant as we assume that  $\{\psi_j\}$  is absolutely summable.

## 13.7 ARIMA Models and Forecasting

If a series  $\{Y_t\}$  is nonstationary in the mean in the sense that there is a trend, then the first difference  $\{X_t = Y_t - Y_{t-1}\}$  or the second difference  $\{X_t = Y_t - 2Y_{t-1} + Y_{t-2}\}$  may behave like a stationary series, and an ARMA model may be used for the differenced series  $\{X_t\}$ . Thus if the  $d$ th-order difference of the sequence  $\{Y_t\}$  follows an  $ARMA(p, q)$  model, then we say that the series  $\{Y_t\}$  follows an  $ARIMA(p, d, q)$  model [ARIMA stands for “integrated autoregressive-moving average”].

Let us briefly discuss forecasting when  $d = 1$ , which is often used in practice. Let  $X_t = Y_t - Y_{t-1}$ . If the observations are  $Y_1, \dots, Y_n$ , then we have the differenced values  $X_2 = Y_2 - Y_1, \dots, X_n = Y_n - Y_{n-1}$ . If  $\{X_t\}$  is modeled by  $ARMA(p, q)$ , we can obtain forecasts  $\hat{X}_{n+h}$ ,  $h = 1, 2, \dots$ . Since

$$Y_{n+h} = X_{n+h} + \cdots + X_{n+1} + Y_n,$$

the forecast of  $Y_{n+h}$  is

$$\hat{Y}_{n+h} = \hat{X}_{n+h} + \cdots + \hat{X}_{n+1} + Y_n.$$

The formula for the mean square error for prediction is a bit more complicated in comparison to the stationary case. Following the arguments in [Section 13.6.4](#),

$$\begin{aligned}\hat{X}_{n+l} - X_{n+l} &= \sum_{j=1}^l \psi_{l-j} \varepsilon_{n+j}, \text{ where } \psi_0 = 1, \text{ and} \\ \hat{Y}_{n+h} - Y_{n+h} &= \sum_{l=1}^h (\hat{X}_{n+l} - X_{n+l}) \\ &= \sum_{l=1}^h \sum_{j=1}^l \psi_{l-j} \varepsilon_{n+j} = \sum_{j=1}^h \left( \sum_{l=0}^{h-j} \psi_l \right) \varepsilon_{n+j},\end{aligned}$$

where  $\{\varepsilon_t\}$  are the error terms in ARMA model for  $\{X_t\}$ . Thus the mean square error for predicting  $Y_{n+h}$  is

$$\sigma^2(h) = \text{E}\left[(\hat{Y}_{n+h} - Y_{n+h})^2 | Y_1, \dots, Y_n\right] = \sigma^2 \sum_{j=1}^h \left( \sum_{l=0}^{h-j} \psi_l \right)^2.$$

Unlike in the stationary case (Section 13.6.4),  $\sigma^2(h)$  here increases linearly with  $h$  as the expression above shows. Thus the prediction interval for  $Y_{n+h}$  can be wide unless  $h$  is small.

## 13.8 Parameter Estimation

We now discuss how to estimate the parameters of ARMA models based on the data  $X_1, \dots, X_n$  and write down the asymptotic distributions of the parameter estimates. For the Gaussian series, the maximum likelihood method can be used to estimate the parameters. However, the actual implementation may always not be easy due to the dependence of the observations and, often appropriate approximations to the likelihood are used in order ease the computation. Details can be found in the book by Box et al. [62]. Here we will basically focus on the least squares type methods. We begin with estimation of the parameters of AR models since it is simpler than the MA or ARMA models. For the discussion below, we assume that the series  $\{X_t\}$  is stationary, invertible, and identifiable.

### 13.8.1 Parameter Estimation: $AR(p)$ Models

If a series  $\{X_t\}$  is  $AR(p)$ , we may estimate the parameters by the method of least squares since it can be rewritten as

$$\begin{aligned} X_t &= \phi_0 + \phi_1 X_{t-1} + \cdots + \phi_p X_{t-p} + \varepsilon_t, \text{ where} \\ \phi_0 &= \mu - (\phi_1 + \cdots + \phi_p)\mu. \end{aligned}$$

Thus one may minimize  $\sum_{t=p+1}^n [X_t - \phi_0 - \phi_1 X_{t-1} - \cdots - \phi_p X_{t-p}]^2$  with respect to  $\phi_0, \dots, \phi_p$  in order to obtain their least squares estimates. However, numerically more stable estimates are obtained via Yule-Walker equations. Recall that the theoretical Yule-Walker equations  $\boldsymbol{\Gamma}_p \boldsymbol{\phi} = \boldsymbol{\gamma}_p$ ,  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_p)^T$ , as given in Eq. (3) are obtained by minimizing  $Q = E[(X_t - \mu) - \phi_1(X_{t-1} - \mu) - \cdots - \phi_p(X_{t-p} - \mu)]^2$  with respect to  $\phi_1, \dots, \phi_p$ . If  $\{\gamma(j)\}$  in Yule-Walker equations are replaced by their empirical estimates  $\{\hat{\gamma}(j)\}$ , then we get the empirical version of the equations  $\hat{\boldsymbol{\Gamma}}_p \boldsymbol{\phi} = \hat{\boldsymbol{\gamma}}_p$ , that is,  $\hat{\boldsymbol{\phi}} = \hat{\boldsymbol{\Gamma}}_p^{-1} \hat{\boldsymbol{\gamma}}_p$ . The Yule-Walker estimate  $\hat{\boldsymbol{\phi}}$  can also be obtained by padding the data  $X_1, \dots, X_n$  as follows. Create  $\tilde{X}_t$ ,  $t = -p+1, \dots, n+p$ , where  $\tilde{X}_t = X_t - \bar{X}$ ,  $t = 1, \dots, n$  and  $\tilde{X}_t = 0$  otherwise. Then one can minimize the least squares criterion

$$\sum_{t=1}^{n+p} (\tilde{X}_t - \phi_1 \tilde{X}_{t-1} - \cdots - \phi_{t-p} \tilde{X}_{t-p})^2$$

with respect to  $\phi_1, \dots, \phi_p$  and the normal equations are the same as the Yule-Walker equations  $\hat{\boldsymbol{\Gamma}}_p \boldsymbol{\phi} = \hat{\boldsymbol{\gamma}}_p$ .

**Theorem 13.8.1.** Let  $\hat{\phi}$  be the solution of the Yule-Walker equations  $\hat{\Gamma}_p \phi = \hat{\gamma}_p$ , where  $\hat{\Gamma}$  and  $\hat{\gamma}_p$  are estimates of  $\Gamma_p$  and  $\gamma_p$  based on the observations  $X_1, \dots, X_n$  from an AR( $p$ ) process. Assuming that  $E[\varepsilon_t^4] < \infty$ ,  $\sqrt{n}(\hat{\phi} - \phi) \xrightarrow{D} N_p(0, \sigma^2 \Gamma_p^{-1})$ .

*Remark 13.8.1.* The above result is also true if the true model is AR( $p$ ) but we fit an autoregressive model of order higher than  $p$ . For  $h > p$ , let  $\phi_h = \begin{pmatrix} \phi \\ \mathbf{0} \end{pmatrix}$ , where  $\phi$  is the vector of autoregressive parameters of the AR( $p$ ) model and  $\mathbf{0}$  is a  $(h-p)$ -dim vector of zeros. Let  $\hat{\phi}_h$  be the solution of the Yule-Walker equations for estimating the parameters of an AR( $h$ ) model, that is,  $\hat{\phi}_h = \hat{\Gamma}_h^{-1} \hat{\gamma}_h$ . Then  $\sqrt{n}(\hat{\phi}_h - \phi_h) \xrightarrow{D} N_h(\mathbf{0}, \sigma^2 \Gamma_h^{-1})$  assuming that  $E[\varepsilon_t^4] < \infty$ .

This result can be used to carry out inferences on  $\phi_1, \dots, \phi_p$  including construction of confidence intervals and deciding if an autoregressive term can be dropped from the model.

An outline of the proof will be given below. One may heuristically guess the asymptotic result since the AR( $p$ ) model can be reexpressed as a Gauss-Markov model and one may use the distributional results of the least squares estimates. However, applying the distributional results of the Gauss-Markov model for the autoregressive case requires justifications due to dependence of the observations. We provide an outline of the main arguments used in the derivation of the asymptotic normality of  $\hat{\phi}$ .

Since  $\{X_t\}$  can be centered by subtracting the mean  $\mu$ , we assume that  $E[X_t] = 0$ . The main idea behind the proof is to decompose  $\hat{\gamma}_p$  so that

$$\hat{\gamma}_p = \hat{\Gamma}_p \phi + n^{-1/2} \delta + O_P(n^{-1}),$$

where  $\delta$  is a  $p$ -dim vector whose  $j$ th element is  $\delta_j = n^{-1/2} \sum_{t=1}^n X_t \varepsilon_{t+j}$ . For the moment assume that  $\delta \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \Gamma_p)$ . Then

$$\begin{aligned} \hat{\Gamma}_p(\hat{\phi} - \phi) &= \hat{\gamma}_p - \hat{\Gamma}_p \phi = [\hat{\Gamma}_p \phi + n^{-1/2} \delta + O_P(n^{-1})] - \hat{\Gamma}_p \phi \\ &= n^{-1/2} \delta + O_P(n^{-1}). \end{aligned}$$

It follows from [Theorem 13.3.4](#) that  $\hat{\gamma}(h) - \gamma(h) \xrightarrow{P} 0$ ,  $h = 0, \dots, p-1$ , and hence  $\hat{\Gamma}_p - \Gamma_p \xrightarrow{P} 0$  and  $\hat{\Gamma}_p^{-1} - \Gamma_p^{-1} \xrightarrow{P} 0$ . Therefore,

$$n^{1/2}(\hat{\phi} - \phi) = \hat{\Gamma}_p^{-1} \delta + O_P(n^{-1/2}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \Gamma_p^{-1}).$$

Let us briefly examine why the asymptotic normality of  $\delta$  is valid. It is fairly easy to check that

$$\begin{aligned} E[\delta_j] &= 0, \\ E[\delta_j \delta_k] &= n^{-1} \sum_{t=1}^n E[X_t X_{t+k-j}] = \gamma(k-j), \quad j \leq k, \text{ and} \\ \text{Cov}[\delta] &= \Gamma_p. \end{aligned}$$

In order to establish the asymptotic normality of  $\delta$ , use the Cramér-Wold device, that is, establish the asymptotic normality of a linear function  $c_1\delta_1 + \dots + c_p\delta_p = \mathbf{c}^T\delta$  of  $\delta$ , where  $\mathbf{c} = (c_1, \dots, c_p)^T$  is a vector of constants. Now, denoting  $S_t = X_t(c_1\varepsilon_{t+1} + \dots + c_p\varepsilon_{t+p})$ , we have

$$\mathbf{c}^T\delta = n^{-1/2} \sum_{t=1}^n X_t(c_1\varepsilon_{t+1} + \dots + c_p\varepsilon_{t+p}) = n^{-1/2} \sum_{t=1}^n S_t$$

Since  $\{X_t\}$  is stationary and has a causal representation,  $X_t = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}$  with  $\psi_0 = 0$  and  $\sum |\psi_j| < \infty$ . For any positive integer  $m$ , let  $X_{m,t} = \sum_{j=0}^m \psi_j \varepsilon_{t-j}$  and  $S_{m,t} = X_{m,t}(c_1\varepsilon_{t+1} + \dots + c_p\varepsilon_{t+p})$ . Since  $\{\varepsilon_t\}$  are iid, the process  $\{S_{m,t}\}$  is  $(m+p)$ -dependent and we can use [Theorem 13.3.2](#) (details omitted) in order to establish asymptotic normality of  $\delta$ .

Let us now see why the decomposition of  $\hat{\gamma}_p$  given above is valid. Since we are assuming that  $E[X_t] = 0$ ,  $X_t = \boldsymbol{\phi}^T \mathbf{X}_{t-1} + \varepsilon_t$  where  $\mathbf{X}_{t-1} = (X_{t-1}, \dots, X_{t-p})^T$ . Using [Lemma 13.3.1](#), the  $j$ th element of  $\hat{\gamma}_p$  is

$$\begin{aligned} \hat{\gamma}(j) &= n^{-1} \sum_{t=1}^n X_t X_{t+j} + O_P(n^{-1}) \\ &= n^{-1} \sum_{t=1}^n X_t (\boldsymbol{\phi}^T \mathbf{X}_{t+j-1} + \varepsilon_{t+j}) + O_P(n^{-1}) \\ &= \sum_{l=1}^p \phi_l n^{-1} \sum_{t=1}^n X_t X_{t+j-l} + n^{-1} \sum_{t=1}^n X_t \varepsilon_{t+j} + O_P(n^{-1}) \\ &= \sum_{l=1}^p \phi_l [\hat{\gamma}(j-l) + O_P(n^{-1})] + n^{-1/2} \delta_j + O_P(n^{-1}) \\ &= \sum_{l=1}^p \phi_l \hat{\gamma}(j-l) + n^{-1/2} \delta_j + O_P(n^{-1}), \quad h = 1, \dots, p. \end{aligned}$$

In the matrix notations, these equations can be written as

$$\hat{\gamma}_p = \hat{\Gamma}_p \boldsymbol{\phi} + n^{-1/2} \delta + O_P(n^{-1}).$$

### 13.8.2 Parameter Estimation: $MA(q)$ Models

Parameter estimation for moving average models is more complicated in comparison to autoregressive models and no closed form solution is available. If the series  $\{X_t\}$  has mean  $\mu$  and we have observations  $X_1, \dots, X_n$ , we can estimate  $\mu$  by the sample mean  $\bar{X}$  and then continue to do the analysis based on  $\{X_t - \bar{X}\}$ . For this reason, in order to make the notations simple, we assume that the series has zero mean. We now describe the least squares and the maximum likelihood methods for estimating  $\theta_1, \dots, \theta_q$ .

Let us follow the ideas used in forecasting  $MA(q)$  models. If  $\boldsymbol{\varepsilon}_0 = (\varepsilon_{-q+1}, \dots, \varepsilon_0)^T$  were known, given  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_q)$ , we may calculate  $\varepsilon_t$ 's starting with

$$\varepsilon_1(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0) = X_1 - (\theta_1 \varepsilon_0 + \dots + \theta_q \varepsilon_{-q+1}),$$

and then iteratively compute, when  $\varepsilon_{-q+1}, \dots, \varepsilon_{t-1}$  are available,  $\varepsilon_t$  as

$$\varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0) = X_t - (\theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}), \quad t = 2, \dots, n.$$

It is important to note that each  $\varepsilon_t$  is a linear combination of  $X_t, \dots, X_1$  and  $\boldsymbol{\varepsilon}_0$ . For the MA(1) case, denoting  $\theta_1$  by  $\theta$ , we have

$$\varepsilon_t(\theta, \boldsymbol{\varepsilon}_0) = X_t - \theta X_{t-1} + \dots + (-\theta)^{t-1} X_1 + (-\theta)^t \varepsilon_0.$$

In the invertible case, that is,  $|\theta| < 1$ , dependence of  $\varepsilon_t(\theta, \boldsymbol{\varepsilon}_0)$  on  $\boldsymbol{\varepsilon}_0$  is negligible when  $t$  is not small since  $(-\theta)^t \rightarrow 0$  rapidly. For the general MA( $q$ ) case, the same is true when it is invertible.

In order to obtain the estimate of  $\boldsymbol{\theta}$ , we may minimize  $\sum_{t=1}^n \varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0)^2$ , with respect to  $\boldsymbol{\theta}$ . This estimate  $\hat{\boldsymbol{\theta}}$  depends on  $\boldsymbol{\varepsilon}_0$ , even though this dependence is negligible in the invertible case as pointed out before. How can one implement this estimation method in practice? The following are among many options:

- (i) take  $\boldsymbol{\varepsilon}_0 = \mathbf{0}$ ,
- (ii) obtain the MLE, and
- (iii) obtain a modified MLE.

The first option is often used and is reasonable in the invertible case. In order to see the second and the third options let us note that the joint pdf of  $\boldsymbol{\varepsilon}_0$  and  $X_1, \dots, X_n$ , under normality (ie,  $\{\varepsilon_t\}$  are iid  $N(0, \sigma^2)$ ), is

$$f(\mathbf{x}, \boldsymbol{\varepsilon}_0, \boldsymbol{\theta}, \sigma) = \left( \sqrt{2\pi}\sigma \right)^{-n-q} \exp \left[ -\frac{1}{2\sigma^2} \sum_{t=1}^n \varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0)^2 - \frac{1}{2\sigma^2} \|\boldsymbol{\varepsilon}_0\|^2 \right],$$

where  $\mathbf{x} = (x_1, \dots, x_n)^T$ . In order to obtain the MLE (option (ii)), one needs to maximize the likelihood after  $\boldsymbol{\varepsilon}_0$  has been integrated out, that is, maximize

$$f(\mathbf{X}, \boldsymbol{\theta}, \sigma) = \int f(\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_0, \sigma) g(\boldsymbol{\varepsilon}_0, \sigma) d\boldsymbol{\varepsilon}_0,$$

where  $g$  is the marginal pdf of  $\boldsymbol{\varepsilon}_0$  and  $\mathbf{X} = (X_1, \dots, X_n)^T$ . Clearly this integration is rather difficult since  $\sum_{t=1}^n \varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0)^2$  does not have a simple expression involving  $\boldsymbol{\varepsilon}_0$ . Another option is to approximate this integral  $f(\mathbf{X}, \boldsymbol{\theta}, \sigma)$  by an average of  $M$  iid copies of  $\boldsymbol{\varepsilon}_0$ , that is, approximate  $f(\mathbf{X}, \boldsymbol{\theta}, \sigma)$  by

$$f_M(\mathbf{X}, \boldsymbol{\theta}, \sigma) = M^{-1} \sum_{j=1}^M f(\mathbf{X}, \boldsymbol{\varepsilon}_{0j}, \boldsymbol{\theta}, \sigma),$$

where  $\boldsymbol{\varepsilon}_{0j}$  are iid as  $\boldsymbol{\varepsilon}_0$ . Mathematically  $f_M(\mathbf{x}, \boldsymbol{\theta}, \sigma)$  converges to  $f(\mathbf{x}, \boldsymbol{\theta}, \sigma)$  in probability as  $M \rightarrow \infty$ . In practice,  $M = n$  should be adequate.

In the third option, one may try to maximize  $f(\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_0, \sigma)$  with respect to  $\boldsymbol{\theta}, \sigma$ , and  $\boldsymbol{\varepsilon}_0$ , an idea used in the derivation of mixed model equations in [Chapter 11](#). Thus one may minimize

$$-2 \log f(\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_0, \sigma) = \frac{1}{\sigma^2} \sum_{t=1}^n \varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0)^2 + \frac{1}{\sigma^2} \|\boldsymbol{\varepsilon}_0\|^2 + (n+q) \log(2\pi\sigma^2),$$

with respect to  $\boldsymbol{\theta}$ ,  $\sigma$ , and  $\boldsymbol{\varepsilon}_0$ . Note that the first two terms in the last expression have the same denominator and hence in order to obtain estimates of  $\boldsymbol{\theta}$  and  $\boldsymbol{\varepsilon}_0$  one needs to minimize the penalized criterion

$$\sum_{t=1}^n \varepsilon_t(\boldsymbol{\theta}, \boldsymbol{\varepsilon}_0)^2 + \|\boldsymbol{\varepsilon}_0\|^2$$

with respect to  $\boldsymbol{\theta}$  and  $\boldsymbol{\varepsilon}_0$ .

For the least squares method, once we have obtained  $\hat{\boldsymbol{\theta}}$ , the estimate of  $\sigma^2$  can be obtained as  $\hat{\sigma}^2 = (n)^{-1} \sum \varepsilon_t(\hat{\boldsymbol{\theta}}, \boldsymbol{\varepsilon}_0)^2$ , where either  $\boldsymbol{\varepsilon}_0 = \mathbf{0}$  or it is estimated by optimizing a penalized criterion as described above.

The following asymptotic result holds for the estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$ . In order to describe the covariance matrix of the asymptotic distribution of  $\hat{\boldsymbol{\theta}}$ , we will adopt some simple notations as in Brockwell and Davis [63]. Let  $\{D_t\}$  be a mean zero  $AR(q)$  process of the form  $\theta(B)D_t = \delta_t$ , where  $\theta(B) = 1 + \theta_1B + \dots + \theta_qB^q$  and  $\{\delta_t\}$  are iid with mean 0 and variance  $\sigma^2$ . Let  $\mathbf{V} = \text{Cov}[D_t]$ , where  $\mathbf{D}_t = [D_1, \dots, D_q]^T$ .

**Theorem 13.8.2.** *Let  $X_1, \dots, X_n$  be observations from a Gaussian  $MA(q)$  series which is invertible. Then as  $n \rightarrow \infty$ ,  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{D} N_q(\mathbf{0}, \sigma^2 \mathbf{V}^{-1})$ , where  $\mathbf{V}$  is as given above and  $\hat{\boldsymbol{\theta}}$  is estimated using any of the methods described above.*

This result makes it possible to obtain the standard errors of the estimates  $\hat{\theta}_1, \dots, \hat{\theta}_q$  and construct confidence intervals for  $\theta_1, \dots, \theta_p$  or carry out tests of hypotheses.

### 13.8.3 Parameter Estimation: $ARMA(p, q)$ Models

As in the  $MA(q)$  case, parameter estimates for ARMA models do not have closed-form expressions. There are many methods for the estimation of parameters and many textbooks on time series provide details of these methods. As in the  $MA(q)$  case, we can center the observations by the sample mean so that we can assume the mean of the series to be zero. Here we outline a simple least squares type method extending some of the ideas outlined for the  $AR(p)$  and  $MA(q)$  cases. Details can be found in the well-known book on time series by Box et al. [62].

We assume that  $\{X_t\}$  is mean zero stationary Gaussian  $ARMA(p, q)$  series which is causal and invertible. If  $\boldsymbol{\varepsilon}_p = (\varepsilon_{p-q+1}, \dots, \varepsilon_p)^T$  were available, then for given values of  $\boldsymbol{\phi}$  and  $\boldsymbol{\theta}$ , we can obtain

$$\varepsilon_{p+1}(\boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_p) = (X_{p+1} - \phi_1X_p - \dots - \phi_pX_1) - (\theta_1\varepsilon_p + \dots + \theta_q\varepsilon_{p-q+1}).$$

Once  $\varepsilon_{p+1}(\boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_p), \dots, \varepsilon_{t-1}(\boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_p)$  are obtained, then it is possible to calculate

$$\varepsilon_t(\boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\varepsilon}_p) = (X_t - \phi_1X_{t-1} - \dots - \phi_pX_{t-p}) - (\theta_1\varepsilon_{t-1} + \dots + \theta_q\varepsilon_{t-q}),$$

$t = p + 2, \dots, n$ . Following the ideas in the  $MA(q)$  case, in order to obtain estimates  $\hat{\phi}$  and  $\hat{\theta}$ , we may minimize  $\sum_{t=p+1}^n \varepsilon_t(\phi, \theta, \epsilon_p)^2$ , with  $\epsilon_p = \mathbf{0}$ , with respect to  $\phi$  and  $\theta$ , or we may minimize

$$\sum_{t=p+1}^n \varepsilon_t(\phi, \theta, \epsilon_p)^2 + \|\epsilon_p\|^2,$$

with respect to  $\phi$ ,  $\theta$ , and  $\epsilon_p$ .

As described in the  $AR(p)$  estimation case, we may also used the idea of padding the data at the beginning and at the end to obtain  $\{\tilde{X}_t: t = -p + 1, \dots, n + p\}$  and then taking  $\epsilon_0 = (\varepsilon_{-q+1}, \dots, \varepsilon_0)^T$ , we can obtain  $\varepsilon_t(\phi, \theta, \epsilon_0)$ ,  $t = 1, \dots, n + p$ , and then obtain the estimates  $\hat{\phi}$  and  $\hat{\theta}$  by minimizing  $\sum_{t=p+1}^n \varepsilon_t(\phi, \theta, \epsilon_0)^2$ , with  $\epsilon_0 = \mathbf{0}$ , or by minimizing the penalized criterion

$$\sum_{t=1}^{n+p} \varepsilon_t(\phi, \theta, \epsilon_0)^2 + \|\epsilon_0\|^2$$

with respect to  $\phi$ ,  $\theta$ , and  $\epsilon_0$ .

An estimate of  $\sigma^2$  is

$$\hat{\sigma}^2 = (n)^{-1} \sum_{t=p+1}^n \varepsilon_t(\hat{\phi}, \hat{\theta}, \hat{\epsilon}_0)^2,$$

where  $\hat{\epsilon}_0$  is  $\mathbf{0}$  or is obtained by optimizing a penalized criterion as described above.

We now write down the asymptotic distribution of  $(\hat{\phi}, \hat{\theta})$ . Let  $\{C_t\}$  and  $\{D_t\}$  be mean zero  $AR(p)$  and  $AR(q)$  series

$$\phi(B)C_t = \delta_t, \quad \theta(B)D_t = \delta_t,$$

where  $\phi(B) = 1 - \phi_1B - \dots - \phi_pB^p$ ,  $\theta(B) = 1 + \theta_1B + \dots + \theta_qB^q$  and  $\{\delta_t\}$  are iid with mean 0 and variance  $\sigma^2$ . Let  $\mathbf{R} = [C_1, \dots, C_p, D_1, \dots, D_q]$  and  $\mathbf{V} = \text{Cov}[\mathbf{R}]$ .

**Theorem 13.8.3.** Let  $X_1, \dots, X_n$  be observations from a stationary Gaussian  $ARMA(p, q)$  series which is invertible and nonredundant. Let  $\beta$  the  $(p+q) \times 1$  vector of  $\phi$  and  $\theta$  stacked vertically and, similarly, let  $\hat{\beta}$  be the stacked vector of  $\hat{\phi}$  and  $\hat{\theta}$  which are obtained using any of the least squares methods outlined above. Then  $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{\mathcal{D}} N_{p+q}(\mathbf{0}, \sigma^2 \mathbf{V}^{-1})$ , where  $\mathbf{V}$  is described above.

Approximate standard errors of the estimates of  $\hat{\phi}$  and  $\hat{\theta}$  can be obtained using this theorem.

## 13.9 Selection of an Appropriate ARMA model

As in any statistical method, model selection is an important part of time series analysis. As in the case of linear models, it is possible to select an appropriate ARMA model by using a criterion such as AIC or BIC. We also introduce another widely used criterion known as AICC which provides a small sample correction to AIC. For an  $ARMA(p, q)$  model, the total number of estimated parameters (excluding  $\sigma^2$ ) is  $k = p + q + 1$ . As in the case of linear

models (Section 11.8.3 in Chapter 11), the  $AIC$ ,  $AICC$ , and  $BIC$  values with  $k = p + q + 1$ , can be written as

$$\begin{aligned} AIC(k) &= -2 \log L(\hat{\mu}_k, \hat{\phi}_k, \hat{\theta}_k, \hat{\sigma}_k) + 2k, \quad BIC(k) = -2 \log L(\hat{\mu}_k, \hat{\phi}_k, \hat{\theta}_k, \hat{\sigma}_k) + [\log(n)]k, \\ AICC(k) &= AIC(k) + 2k(k+1)/(n-k-1), \end{aligned}$$

where  $L$  is the likelihood, and  $(\hat{\mu}_k, \hat{\phi}_k, \hat{\theta}_k, \hat{\sigma}_k)$  is the MLE of  $(\mu, \phi, \theta, \sigma)$  under the  $ARMA(p, q)$  model. Statistical computing packages such as R or MATLAB can calculate the values of these criteria for an ARMA model given a data set. As usual, an appropriate model is selected by minimizing the criterion of choice.

## 13.10 Spectral Analysis

Analysis of time series has two important aspects

- (a)** model fitting and forecasting, and
- (b)** understanding of the hidden periodicities.

We have discussed the first aspect in the previous sections, and we now discuss the second which involves a deeper analysis of the spectral density function defined in Eq. (1).

The basis of spectral analysis is an important mathematical result which states that any stationary time series (causal or noncausal) can be approximated by linear combinations of sines and cosines with random coefficients. Toward this end, let us define an important quantity: frequency. A frequency is a real number between 0 and 1/2. For the monthly energy (electricity) data mentioned at the beginning of this chapter, once the trend is estimated and subtracted from the data, the detrended series has a similar pattern of behavior every 12 time points (months). Energy consumptions in January are similar, energy consumptions in March are similar, and so on. In such a case, we can say that energy consumption has an “important” frequency at  $w = 1/12$ . A plot of the annual sunspots recorded over the last  $n = 313$  years reveals that the peaks are occurring between 8 and 12 years. However, unlike in the energy data (which is seasonal), the times of the peaks are not fixed in the sunspot series. There are 28.5 cycles in the series. So the series tends to behave similarly every  $313/28.5 = 10.98$  years (on the average). Thus there is a peak at frequency at  $w = 1/10.98 = 0.091$ .

Consider a series  $X_t = A \cos(2\pi wt) + B \sin(2\pi wt)$ ,  $0 \leq w \leq 1/2$ , where  $A$  and  $B$  are mutually uncorrelated rv's with mean 0 and variance  $\sigma^2$ . Then  $E[X_t] = 0$ ,  $\text{Var}[X_t] = \sigma^2$ , and the series  $\{X_t\}$  is stationary since  $\text{Cov}[X_t, X_{t+h}] = \sigma^2 \cos(2\pi wh)$  depends only on  $h$ . Note that  $\text{Corr}[X_t, X_{t+h}] = \cos(2\pi hw) = 1$  whenever  $hw$  is a positive integer. Thus if  $1/w$  is an integer, then the series repeats itself at every  $1/w$  time points, that is,  $X_t = X_{t+h}$  when  $h$  is an integer multiple of  $1/w$ . However, even if  $1/w$  is not a rational number, the correlation  $\cos(2\pi wh)$  is high whenever  $hw$  is close to an integer. We call this series an elementary periodic series with frequency  $w$  and variance  $\sigma^2$ .

Now consider a series  $X_t$  which is a sum of  $M$  such elementary periodic series with distinct frequencies  $w_1, \dots, w_M$  and variances  $\sigma_1^2, \dots, \sigma_M^2$ , that is,

$$X_t = \sum_{j=1}^M \{A_j \cos(2\pi w_j t) + B_j \sin(2\pi w_j t)\}, \quad (7)$$

where  $\{A_j\}$ ,  $\{B_j\}$  are mean zero mutually uncorrelated rv's with  $\text{Var}[A_j] = \text{Var}[B_j] = \sigma_j^2$ . It is also assumed that  $\{A_j\}$  and  $\{B_j\}$  are also uncorrelated to each other. Then it can be seen that  $\{X_t\}$  is stationary with

$$\begin{aligned} E[X_t] &= 0, \text{Var}[X_t] = \sigma_1^2 + \dots + \sigma_M^2, \text{ and} \\ \text{Cov}[X_t, X_{t+h}] &= \sum_{j=1}^M \sigma_j^2 \cos(2\pi w_j h) = \gamma(h). \end{aligned} \quad (8)$$

### Remarks

- (a) If  $X_t$  has mean  $\mu$ , then the representation above is valid with  $X_t - \mu$  on the left-hand side of Eq. (7). From now on we assume the mean to be equal to zero since we can always carry out spectral analysis after subtracting the mean from the series.
- (b) For the monthly electricity consumption data  $\{Y_t = \log(\text{sales})\}$ , we have briefly discussed the model  $Y_t = m_t + S_t + X_t$ , where  $\{X_t\}$  is stationary (Section 13.1). There are many approaches to the analysis of such a data set. We may estimate the trend  $m_t$  and the seasonal effect  $S_t$ , and then subtract them from  $Y_t$  in order to get an estimate of the stationary part  $\{X_t\}$ . However, there is another way of modeling this. We can subtract the trend only and the remainder, that is,  $S_t + X_t$  can be often considered stationary, especially if  $\{S_t\}$  is deemed to be stochastic. Thus if a sequence has no trend, but has seasonality whose variance does not depend on time  $t$ , then the sequence itself can be considered stationary.
- (c) In some trivial cases,  $M$  may be small. But in general  $M$  is large.
- (d) The goal of spectral analysis is to find  $\sigma_j$ 's. Since  $\text{Var}[X_t] = \sigma_1^2 + \dots + \sigma_M^2$ , the contribution of the  $j$ th elementary periodic series to this variance (at frequency  $w_j$ ) is  $\sigma_j^2$ . It is of interest to find out which frequencies contribute more to this variability than the others.
- (e) If  $M$  is large (mathematically  $M \rightarrow \infty$ ) and almost all  $\sigma_j$ 's are not equal to zero, then these  $\sigma_j$ 's need to be small so that  $\text{Var}[X_t] = \sigma_1^2 + \dots + \sigma_M^2$  remains finite as  $M \rightarrow \infty$ . This can be done if  $\sigma_j^2 = O(1/M)$ . In Section 13.2.1, it is written that if  $\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty$ , then the spectral density is given by  $f(w) = \sum_{h=-\infty}^{\infty} \gamma(h) \exp(-2\pi i h w)$ . However, for the series in Eq. (7),  $\sum_{h=-\infty}^{\infty} |\gamma(h)|$  is not finite and hence we may consider a truncated version  $f_M(w) = \sum_{h=-M}^M \gamma(h) \exp(-2\pi i h w)$ . It is left as an exercise for the reader to verify that

$$f_M(w) = [(2M+1)/2]\sigma_j^2 \quad \text{if } w = j/(2M+1), \quad j = 1, \dots, M.$$

In order for  $f_M$  to be approximately equal to a bounded spectral density  $f$  for large  $M$ , we need to have  $\sigma_j^2 \approx [2/(2M+1)]f(w_j)$  when  $w_j = j/(2M+1)$  and one of the goals of the spectral analysis is to obtain an estimate of the spectral density  $f$  based on the available data  $X_1, \dots, X_n$ .

### 13.10.1 Representation of a Stationary Series

If  $\{X_t\}$  is a mean zero stationary series, can one approximate it as in Eq. (7)? The answer is yes if the underlying spectral density function  $f$  as defined in Eq. (1) is square integrable. The arguments given here are heuristic and detailed proofs can be found in Gikhman and Skorokhod [64]. Consider the following random functions

$$Z_{1,M}(w) = \sum_{s=-M}^M \frac{\sin(2\pi sw)}{2\pi s} X_s \quad \text{and} \quad Z_{2,M}(w) = \sum_{s=-M}^M \frac{1 - \cos(2\pi sw)}{2\pi s} X_s,$$

with the understanding that when  $s = 0$ ,  $\sin(2\pi sw)/(2\pi s) = w$  and  $[1 - \cos(2\pi sw)]/(2\pi s) = 0$ . Then it is easy to check that for  $-M \leq t \leq M$ ,

$$X_t = 2 \int_0^{1/2} \cos(2\pi tw) dZ_{1,M}(w) + 2 \int_0^{1/2} \sin(2\pi tw) dZ_{2,M}(w).$$

Now consider the limiting random functions  $Z_1$  of  $Z_{1,M}$  and  $Z_2$  of  $Z_{2,M}$  as  $M \rightarrow \infty$

$$Z_1(w) = \sum_{s=-\infty}^{\infty} \frac{\sin(2\pi sw)}{2\pi s} X_s \quad \text{and} \quad Z_2(w) = \sum_{s=-\infty}^{\infty} \frac{1 - \cos(2\pi sw)}{2\pi s} X_s,$$

which exist for  $0 \leq w < 1/2$  in the sense that  $E[Z_1(w)^2] < \infty$  and  $E[Z_2(w)^2] < \infty$ . Let us denote  $R_{t,M} = \int_0^{1/2} \cos(2\pi tw) dZ_{1,M}(w)$ . For any given  $t$ , the sequence  $\{R_{t,N}\}$  is Cauchy in the mean square sense, that is,  $E[(R_{t,M} - R_{t,N})^2] \rightarrow 0$  and  $M, N \rightarrow \infty$ . Thus a limit of  $R_{t,M}$  exists in the mean square sense as  $M \rightarrow \infty$  and the limit is denoted by

$$\int_0^{1/2} \cos(2\pi tw) dZ_1(w).$$

A similar argument can be used to show that the limit of  $\int_0^{1/2} \sin(2\pi tw) dZ_{2,M}(w)$  exists in the mean square sense as  $M \rightarrow \infty$  and the limit is denoted by

$$\int_0^{1/2} \sin(2\pi tw) dZ_2(w).$$

Thus we can represent the time series  $\{X_t\}$  as

$$X_t = 2 \int_0^{1/2} \cos(2\pi tw) dZ_1(w) + 2 \int_0^{1/2} \sin(2\pi tw) dZ_2(w). \quad (9)$$

The random function  $Z_1$  has the *orthogonal increment* property, that is, if  $w_1 \neq w_2$ , then  $Z_1(w_1 + \delta) - Z_1(w_1)$  and  $Z_1(w_2 + \delta) - Z_1(w_2)$ ,  $\delta > 0$ , are uncorrelated if the intervals  $(w_1, w_1 + \delta]$  and  $(w_2, w_2 + \delta]$  are disjoint. Moreover,  $E[\{Z_1(w + \delta) - Z_1(w)\}^2] \approx (\delta/2)f(w)$  when  $\delta > 0$  is small enough. The random function  $Z_2$  also has the same

property as  $Z_1$ . Additionally,  $Z_1(w_1)$  is uncorrelated with  $Z_2(w_2)$  for any  $w_1, w_2$ . The representation of  $X_t$  as given above in Eq. (9) is known as the *Cramér Representation* of a stationary series.

Now consider approximating the integral  $2 \int_0^{1/2} \cos(2\pi tw) dZ_1(w)$  by a finite sum as follows. For a positive integer  $M$ , denote  $w_j = j/(2M+1)$ ,  $I_j = (w_{j-1}, w_j]$ ,  $Z_1(I_j) = Z_1(w_j) - Z_1(w_{j-1})$ ,  $j = 1, \dots, M$ . Ignoring the integral over  $(w_M, 1/2)$ , which is reasonable when  $M$  is large, we have

$$\begin{aligned} 2 \int_0^{1/2} \cos(2\pi tw) dZ_1(w) &\approx 2 \sum_{j=1}^M \int_{I_j} \cos(2\pi tw) dZ_1(w) \\ &\approx 2 \sum_{j=1}^M \cos(2\pi tw_j) Z_1(I_j) = \sum_{j=1}^M A_j \cos(2\pi tw_j), \end{aligned}$$

with  $A_j = 2Z_1(I_j)$ . Due to the orthogonal increment property of the random function  $Z_1$ , the rv's  $\{A_j\}$  are mutually uncorrelated and

$$\text{Var}[A_j] = 4\text{Var}[Z_1(I_j)] \approx 4[1/(2(2M+1))]f(w_j) = [2/(2M+1)]f(w_j).$$

Similarly we can approximate

$$2 \int_0^{1/2} \sin(2\pi tw) dZ_2(w) \approx \sum_{j=1}^M B_j \cos(2\pi tw_j),$$

where  $B_j = 2Z_2(I_j)$ . Here  $\{B_j\}$  are also mutually uncorrelated with  $\text{Var}[B_j] \approx [2/(2M+1)]f(w_j)$ . Moreover,  $\{A_j\}$  and  $\{B_j\}$  are also uncorrelated with each other. Thus combining all the arguments above we have

$$X_t \approx \sum_{j=1}^M \{A_j \cos(2\pi w_j t) + B_j \sin(2\pi w_j t)\},$$

when  $M$  is large.

*Remark 13.10.1.* The Cramér representation of  $X_t$  as given in Eq. (9) is true more generally under weaker conditions than given here (Chapter 5 in Gikhman and Skorokhod [64]), but we will not concern ourselves with such mathematical details.

### 13.10.2 Periodogram

Consider the series in Eq. (7) with  $w_j = j/n$ ,  $j = 1, \dots, M$ , where  $M$  is the largest integer for which  $M/n < 1/2$ . We can estimate  $A_j$  and  $B_j$  from the data  $\{X_1, \dots, X_n\}$  using the method of least squares

$$\begin{aligned} \hat{A}_j &= \frac{\sum_{t=1}^n X_t \cos(2\pi w_j t)}{\sum_{t=1}^n \cos^2(2\pi w_j t)}, \quad \hat{B}_j = \frac{\sum_{t=1}^n X_t \sin(2\pi w_j t)}{\sum_{t=1}^n \sin^2(2\pi w_j t)}, \text{ or} \\ \hat{A}_j &= (2/n) \sum_{t=1}^n X_t \cos(2\pi w_j t), \quad \hat{B}_j = (2/n) \sum_{t=1}^n X_t \sin(2\pi w_j t). \end{aligned}$$

The last equalities hold since it can be shown that

$$\sum_{t=1}^n \cos^2(2\pi w_j t) = n/2 \quad \text{and} \quad \sum_{t=1}^n \sin^2(2\pi w_j t) = n/2.$$

The quantity  $P(w_j) = \hat{A}_j^2 + \hat{B}_j^2$  is called the scaled periodogram and its rescaled version

$$I(w_j) = (n/4)P(w_j)$$

is called the periodogram. The main use of the periodogram  $I(w_j)$  is as an estimator of  $f(w_j)$ , the spectral density function at frequency  $w_j = j/n$ . There is a related quantity called the discrete Fourier transform of the data

$$\begin{aligned} d(w_j) &= n^{-1/2} \sum_{t=1}^n X_t \exp(-2\pi i w_j t) \\ &= n^{-1/2} \sum_{t=1}^n X_t \cos(2\pi w_j t) - i n^{-1/2} \sum_{t=1}^n X_t \sin(2\pi w_j t) \\ &:= \tilde{X}_{c,n}(w_j) - i \tilde{X}_{s,n}(w_j) \end{aligned}$$

where  $i = \sqrt{-1}$  is the imaginary number, and  $\tilde{X}_{c,n}(w_j)$  and  $\tilde{X}_{s,n}(w_j)$  are the discrete cosine and sine transforms introduced in [Section 13.2.1](#). The connection between the discrete Fourier transform and the periodogram is

$$I(w_j) = |d(w_j)|^2 = \tilde{X}_{c,n}^2 + \tilde{X}_{s,n}^2.$$

In the general case when the series  $\{X_t\}$  may not have zero mean, the definitions of  $\tilde{X}_{c,n}$ ,  $\tilde{X}_{s,n}$ ,  $d$ , and  $I$  are based on centered observations  $\{X_t - \bar{X}\}$ , that is,

$$\begin{aligned} \tilde{X}_{c,n}(w_j) &= n^{-1/2} \sum_{t=1}^n (X_t - \bar{X}) \cos(2\pi w_j t), \\ \tilde{X}_{s,n}(w_j) &= n^{-1/2} \sum_{t=1}^n (X_t - \bar{X}) \sin(2\pi w_j t), \\ d(w_j) &= \tilde{X}_{c,n}(w_j) - i \tilde{X}_{s,n}(w_j), \quad I(w_j) = |d(w_j)|^2 = \tilde{X}_{c,n}^2 + \tilde{X}_{s,n}^2, \end{aligned}$$

and the periodogram values are obtained at frequencies  $w_j = j/n$ ,  $j = 1, \dots, M$ . We should note that the periodogram values at frequencies  $\{w_j\}$  as defined in the last displayed identities can be, and are usually, calculated when observations  $X_1, \dots, X_n$  from a stationary series are available.

If  $\{X_t\}$  is a mean zero Gaussian series, then clearly  $\tilde{X}_{c,n}(w_j)$  and  $\tilde{X}_{s,n}(w_j)$  are normally distributed since each of them is a linear function of  $X_1, \dots, X_n$ . However, asymptotic normality of  $\tilde{X}_{c,n}(w_j)$  and  $\tilde{X}_{s,n}(w_j)$  hold if the observations are from a linear process as given in Eq. (2) using arguments outlined in the proof of [Theorem 13.3.1](#).

**Theorem 13.10.1.** Let  $\{X_t\}$  be a mean zero linear process as given in Eq. (2) and assume that  $0 < w_j = j/n < 1/2$ ,  $j = 1, \dots, M$ . As  $n \rightarrow \infty$ ,

- (a)  $\tilde{X}_n(w_j) = (\tilde{X}_{c,n}(w_j), \tilde{X}_{s,n}(w_j))^T \xrightarrow{\mathcal{D}} N_2(\mathbf{0}, (1/2)f(w_j)\mathbf{I})$ ,  $j = 1, \dots, M$ , where  $\mathbf{I}$  is the  $2 \times 2$  identity matrix, and  $\tilde{X}_n(w_1), \dots, \tilde{X}_n(w_M)$  are asymptotically independent,
- (b)  $I(w_j)$  is approximately distributed as  $(\xi_j/2)f(w_j)$ , where  $\xi_j \sim \chi_2^2$ ,  $j = 1, \dots, M$ , and  $I(w_1), \dots, I(w_M)$  are approximately independent.

### 13.10.3 Estimation of the Spectral Density

Since  $I(w_j)$  is approximately distributed as  $(\xi_j/2)f(w_j)$ , where  $\xi_j \sim \chi_2^2$ , and  $E[I(w_j)]$  is approximately equal to  $f(w_j)$  (noting that  $E[\xi_j] = 2$ ). Thus  $I(w_j)$  is almost an unbiased estimator of  $f(w_j)$ . However, it is not a consistent estimator since  $\text{Var}[I(w_j)]$  is approximately equal to  $\text{Var}[\xi_j/2]f(w_j)^2 = f(w_j)^2$  and this variance does not converge to zero as the sample size  $n \rightarrow \infty$ . In order to construct a consistent estimator of the spectral density at any point  $w \in (0, 1/2)$ , we may use a weighted average of  $I(w_j)$  for  $w_j$  in a small neighborhood of  $w$ . Assuming that the spectral density function  $f$  is smooth, such a weighted average should lead to a better estimate of  $f(w)$  since  $I(w_j)$ 's are asymptotically independent. Toward this purpose, we may use the kernel method for regression discussed in [Chapter 9](#). Let  $K$  be a pdf on  $[-1, 1]$  and, as in [Chapter 9](#), assume that

$$(i) K \text{ is symmetric about } 0, \quad (ii) \int zK(z)dz = 0.$$

In the discussion below it is assumed that periodogram values at frequencies 0 and  $1/2$  are excluded, so the total number of frequencies  $M$  may be smaller than  $[n/2]$ .

We can now obtain a nonparametric estimate of  $f(w)$ ,  $w \in (0, 1/2)$ , with kernel  $K$  and bandwidth  $h_n \rightarrow 0$  as

$$\hat{f}(w) = \frac{\sum_{j=1}^M K(h_n^{-1}(w - w_j))I(w_j)}{\sum_{j=1}^M K(h_n^{-1}(w - w_j))}.$$

The bias and variance properties of the kernel estimate  $\hat{f}(w)$  are similar to those in [Chapter 9](#). Writing  $K_j = K(h_n^{-1}(w - w_j))$  and  $K_+ = \sum K_j$ , we have

$$\begin{aligned} E[\hat{f}(w)] - f(w) &\approx \sum K_j E[\xi_j/2]f(w_j)/K_+ - f(w) \\ &= \sum K_j f(w_j)/K_+ - f(w) = \sum K_j [f(w_j) - f(w)]/K_+. \end{aligned}$$

If  $f$  is twice differentiable and  $f''$  is continuous, then a two term Taylor expansion yields

$$f(w_j) - f(w) = (w_j - w)f'(w) + (1/2)(w_j - w)^2f''(w'_j),$$

where  $w'_j$  lies between  $w$  and  $w_j$ . Employing arguments similar to the ones used in [Chapter 9](#), we get the bias of  $\hat{f}(w)$  as

$$\begin{aligned}\mathbb{E}[\hat{f}(w)] - f(w) &\approx \sum K_j [f(w_j) - f(w)] / K_+ \\ &= (1/2) h_n^{-2} f''(w) \int z^2 K(z) dz [1 + o(1)] \\ &= (1/2) h_n^2 f''(w) \mu_2(K) [1 + o(1)],\end{aligned}$$

where  $\mu_2(K) = \int z^2 K(z) dz$ .

Approximate independence of  $I(w_j)$ 's may be used to calculate the variance of  $\hat{f}(w)$  and thus we have

$$\begin{aligned}\text{Var}[\hat{f}(w)] &\approx \sum K_j^2 \text{Var}[\xi_j/2] f(w_j)^2 / K_+^2 = \sum K_j^2 f(w_j)^2 / K_+^2 \\ &\approx \left[ \sum K_j^2 / K_+^2 \right] f(w)^2,\end{aligned}\tag{10}$$

where the last step is justified since  $f(w_j)$  is approximately equal to  $f(w)$  by continuity when  $|w_j - w| \leq h_n$ . Once again, employing arguments similar to the ones used in Chapter 9, we have

$$\text{Var}[\hat{f}(w)] \approx (nh_n)^{-1} f(w)^2 \int K^2(z) dz = (Mh_n)^{-1} f(w)^2 \|K\|^2,$$

where  $\|K\|^2 = \int K^2(z) dz$ .

Since the mean square error of any estimator is the sum of its variance and square of its bias, the mean square error of  $\hat{f}(w)$  is approximately given by

$$(nh_n)^{-1} f(w)^2 \|K\|^2 + (1/4) h_n^4 f''(w)^2 \mu_2(K)^2.$$

The last expression is convex in  $h_n$  and it is minimized at  $h_n^* = c_1 n^{-1/5}$ , where  $c_1 = [f(w) \|K\| / \{f''(w) \mu_2(K)\}]^{2/5}$ . The minimum mean square error of  $\hat{f}(w)$  (at  $h_n = h_n^*$ ) is approximately equal to  $c_2 n^{-4/5}$ , where  $c_2$  is a constant that depends  $f(w)$ ,  $f''(w)$ ,  $\|K\|$ , and  $\mu_2(K)$ .

*Remark 13.10.2.*

- (a) It is possible to obtain asymptotic normality of  $\hat{f}(w)$  as given in Chapter 9 and the results are similar.
- (b) As discussed in Section 9.6 of Chapter 9, a drawback of the kernel density or kernel regression estimates is that near the boundary points of the independent variable, the bias may be of order  $h_n$  and not  $h_n^2$ . However, in the case of spectral density estimation, this does not pose a problem since the spectral density  $f$  is periodic and symmetric about 0, so one can obtain a periodogram estimate at point  $-w_j$ , where  $0 < w_j < 1/2$ , by taking  $I(-w_j) = I(w_j)$ . Even though  $I(0) = 0$  (follows from the formula), we can obtain an estimate  $\hat{f}(0)$  of  $f(0)$  using the kernel method since the periodogram values at negative frequencies can be obtained as mentioned above. For a frequency near zero, say at  $w = h_n/2$ , all the values of  $I(w_j)$ ,  $|w_j - w| \leq h_n$ , are now available (substituting  $\hat{f}(0)$  for the periodogram at frequency 0) and we can obtain an estimate  $\hat{f}(w)$  of  $f(w)$  using the kernel method. Thus the bias of  $\hat{f}(w)$  is of order  $h_n^2$  when  $w$  is close to zero. However, it should be pointed out that when  $w$  is close to 0,

the variance of  $\hat{f}(w)$  is different from the formula given in Eq. (10) and it needs to be recalculated since  $\hat{f}(w)$  is no longer a weighted average of approximately independent rv's.

Similar strategies can be used to estimate  $f(1/2)$  and  $f(w)$  when  $w$  is close to 1/2.

Since  $f$  is periodic on  $[-1/2, 1/2]$  and symmetric about 0, for any  $1/2 < w < 1$ ,  $f(w) = f(w - 1) = f(1 - w)$  and the same is also true for the periodogram values. For any  $n/2 < j < n$ , we can get  $I(j/n) = I(1 - j/n)$ . Thus a kernel estimate of  $f(w)$  for  $w$  near 1/2 does not have any inadequacy in terms of inflated bias, but its variance needs to be recalculated as it is not the same as the formula given in Eq. (10).

- (c) As in any nonparametric method, one needs to obtain an estimate of the bandwidth in a data dependent manner. One may apply the method of cross-validation as outlined in [Chapters 9](#) and [11](#) for this purpose. It should be pointed out that in the context of kernel regression discussed in [Chapter 9](#), the theoretical justification of the method of cross-validation relies on the assumption that the data consist of iid observations  $(Y_j, X_j), j = 1, \dots, n$ . A kernel estimate of the spectral density  $f(w)$  is based on the data  $\{(I(w_j), w_j), j = 1, \dots, M\}$ , which are not iid. Even though  $w_j$ 's are nonrandom, the use of cross-validation does not pose a problem since  $w_j$ 's are equally spaced.

### 13.10.4 Linear Filtering

For a series  $\{X_t\}$ , it is sometimes of interest to study the behavior of the first difference or a running weighted average of the series such as

- (a)  $Z_t = X_t - X_{t-1}$ ,
- (b)  $Z_t = (1/2)X_t + (1/2)X_{t-1}$ , and
- (c)  $Z_t = (1/3)X_t + (1/3)X_{t-1} + (1/3)X_{t-2}$ .

For each of the three cases above,  $\{Z_t\}$  is a linear combination of  $\{X_t\}$ . A linear combination of  $\{X_t\}$  is called a filtered series of  $\{X_t\}$ . It turns out that there is a nice formula connecting the spectral density of the original series to that of the filtered series when  $\{X_t\}$  is stationary. Denoting the spectral density functions of  $\{X_t\}$  and  $\{Z_t\}$  by  $f_X$  and  $f_Z$ , respectively, the spectral density functions for  $\{Z_t\}$  for (a) and (b) are (justifications given below)

- (a)  $f_Z(w) = [2 - 2 \cos(2\pi w)] f_X(w)$ , and
- (b)  $f_Z(w) = (1/2)[1 + \cos(2\pi w)] f_X(w)$ .

In each case, the spectral density of  $Z_t$  is equal to the spectral density of  $X_t$  times a weight function. Note that in (a), the weight function is zero at  $w = 0$  and it monotonically increases to the value of 4 at  $w = 1/2$ . In other words, the higher the frequency, the higher is the weight, indicating that the first difference of  $\{X_t\}$  is a rougher series than  $\{X_t\}$ . For the second case, the weight function  $(1/2)[1 + \cos(2\pi w)]$  equals 1 at  $w = 0$  and then

it decreases to zero at  $w = 1/2$ . This indicates that the running average  $\{Z_t\}$  of  $\{X_t\}$  is smoother than  $\{X_t\}$  since, for this case, higher frequencies have lower weights.

Consider the filtered series of a mean zero stationary series  $\{X_t\}$ ,

$$Z_t = \sum_{j=-\infty}^{\infty} \psi_j X_{t-j},$$

where  $\{\psi_j\}$  are constants satisfying the condition  $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$ . If  $\gamma_X$  is the autocovariance function of the series  $\{X_t\}$ , then

$$\begin{aligned} E[Z_t^2] &= \sum_{-\infty < j, k < \infty} \psi_j \psi_k \text{Cov}[X_{t-j}, X_{t-k}] \\ &= \sum_{-\infty < j, k < \infty} \psi_j \psi_k \gamma_X(j - k) < \infty, \end{aligned}$$

since  $\gamma_X$  is bounded and  $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$ . The fact that  $\{Z_t\}$  is stationary follows from the fact that

$$\begin{aligned} \text{Cov}[Z_t, Z_{t+h}] &= \sum_{-\infty < j, k < \infty} \psi_j \psi_k \text{Cov}[X_{t-j}, X_{t+h-k}] \\ &= \sum_{-\infty < j, k < \infty} \psi_j \psi_k \gamma_X(h + j - k) \end{aligned}$$

depends only on  $h$ . Thus we conclude that  $\{Z_t\}$  is stationary with the autocovariance function

$$\gamma_Z(h) = \sum_{-\infty < j, k < \infty} \psi_j \psi_k \gamma_X(h + j - k). \quad (11)$$

Now let us look at the spectral density function of the series  $\{Z_t\}$  which is

$$\begin{aligned} f_Z(w) &= \sum_{h=-\infty}^{\infty} \exp(-2\pi i h w) \gamma_Z(h) \\ &= \sum_{h=-\infty}^{\infty} \exp(-2\pi i h w) \sum_{-\infty < j, k < \infty} \psi_j \psi_k \gamma_X(h + j - k) \\ &= \sum_{h=-\infty}^{\infty} \sum_{-\infty < j, k < \infty} [\psi_j \exp(2\pi i j w)][\psi_k \exp(-2\pi i k w)] \\ &\quad \times [\exp(-2\pi i(h + j - k)w) \gamma_X(h + j - k)] \\ &= \sum_{-\infty < j, k < \infty} [\psi_j \exp(2\pi i j w)][\psi_k \exp(-2\pi i k w)] \\ &\quad \times \sum_{h=-\infty}^{\infty} [\exp(-2\pi i(h + j - k)w) \gamma_X(h + j - k)]. \end{aligned}$$

Writing  $l = h + j - k$ , we see that

$$\begin{aligned} & \sum_{h=-\infty}^{\infty} [\exp(-2\pi i(h+j-k)w) \gamma_X(h+j-k)] \\ &= \sum_{l=-\infty}^{\infty} \exp(-2\pi ilw) \gamma_X(l) = f_X(w). \end{aligned}$$

Hence

$$\begin{aligned} f_Z(w) &= \sum_{-\infty < j, k < \infty} [\psi_j \exp(2\pi ijw)][\psi_k \exp(-2\pi ikw)] f_X(w) \\ &= \left| \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi ijw) \right|^2 f_X(w) \\ &= |\Psi(w)|^2 f_X(w), \text{ where} \\ \Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi iwj) \end{aligned} \tag{12}$$

is called the frequency response function. Thus we arrive at the following important result.

**Lemma 13.10.1.** *Let  $\{X_t\}$  be a mean zero stationary series with autocovariance function  $\gamma_X$  and spectral density function  $f_X$ . Consider the filtered series  $Z_t = \sum_{j=-\infty}^{\infty} \psi_j X_{t-j}$  where  $\{\psi_j\}$  are constants satisfying the condition  $\sum_{j=-\infty}^{\infty} |\psi_j| < \infty$ .*

- (a) *The series  $\{Z_t\}$  is mean zero stationary with the autocovariance function  $\{\gamma_Z(h)\}$  given in Eq. (11).*
- (b) *The spectral density function of  $\{Z_t\}$  is  $f_Z(w) = |\Psi(w)|^2 f_X(w)$ , where  $\Psi(w)$  is the frequency response function given in Eq. (12).*

**Example 13.10.1.** Let  $Z_t = X_t - X_{t-1}$ , where  $\{X_t\}$  is mean zero stationary. In this case  $\psi_0 = 1$ ,  $\psi_1 = -1$ , and  $\psi_j = 0$  when  $j \neq 0$  and  $j \neq 1$ . Then the frequency response function is

$$\begin{aligned} \Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi iwj) \\ &= (1) \exp(-2\pi iw0) + (-1) \exp(-2\pi iw) \\ &= 1 - \exp(-2\pi iw) = 1 - \cos(2\pi w) + i \sin(2\pi w), \text{ and} \\ |\Psi(w)|^2 &= [1 - \cos(2\pi w)]^2 + \sin^2(2\pi w) = 2 - 2 \cos(2\pi w). \end{aligned}$$

Hence the spectral density of  $\{Z_t\}$  is

$$f_Z(w) = [2 - 2 \cos(2\pi w)] f_X(w).$$

**Example 13.10.2.** Let  $Z_t = (1/2)X_t + (1/2)X_{t-1}$ , where  $\{X_t\}$  is mean zero stationary. Then  $\psi_0 = \psi_1 = 1/2$  and  $\psi_j = 0$  otherwise. The frequency response function is

$$\begin{aligned}\Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi i w j) \\ &= (1/2) \exp(-2\pi i w 0) + (1/2) \exp(-2\pi i w) \\ &= (1/2)[1 + \exp(-2\pi i w)] \\ &= (1/2)[1 + \cos(2\pi w) - i \sin(2\pi w)], \text{ and} \\ |\Psi(w)|^2 &= (1/4)[1 + \cos(2\pi w)]^2 + (1/4) \sin^2(2\pi w) \\ &= (1/4)[2 + 2 \cos(2\pi w)] = (1/2)[1 + \cos(2\pi w)].\end{aligned}$$

Thus the spectral density function of  $\{Z_t\}$  is

$$f_Z(w) = (1/2)[1 + \cos(2\pi w)]f_X(w).$$

**Example 13.10.3.** Let  $Z_t = (X_t + \dots + X_{t-L+1})/L$ , where  $\{X_t\}$  is mean zero and stationary. Then  $\psi_j = 1/L$  when  $j = 0, \dots, L-1$ , and  $= 0$  otherwise. The frequency response function is

$$\begin{aligned}\Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi i w j) = (1/L) \sum_{j=0}^{L-1} \exp(-2\pi i w j) \\ &= (1/L) \frac{1 - \exp(-2\pi i w L)}{1 - \exp(-2\pi i w)}.\end{aligned}$$

Since

$$\begin{aligned}|1 - \exp(-2\pi i w L)|^2 &= |1 - \cos(2\pi w L) + i \sin(2\pi w L)|^2 \\ &= |1 - \cos(2\pi w L)|^2 + |\sin(2\pi w L)|^2 \\ &= 2 - 2 \cos(2\pi w L), \text{ and similarly} \\ |1 - \exp(-2\pi i w)|^2 &= 2 - 2 \cos(2\pi w),\end{aligned}$$

we have

$$\begin{aligned}|\Psi(w)|^2 &= (1/L^2) \frac{|1 - \exp(-2\pi i w L)|^2}{|1 - \exp(-2\pi i w)|^2} \\ &= (1/L^2) \frac{2 - 2 \cos(2\pi w L)}{2 - 2 \cos(2\pi w)} = (1/L^2) \frac{1 - \cos(2\pi w L)}{1 - \cos(2\pi w)}.\end{aligned}$$

The spectral density function of  $\{Z_t\}$  is

$$f_Z(w) = (1/L^2) \frac{1 - \cos(2\pi w L)}{1 - \cos(2\pi w)} f_X(w).$$

### 13.10.5 Spectral Density for ARMA

In this section we use [Lemma 13.10.1](#) to obtain an explicit expression of the spectral density function of an  $ARMA(p, q)$  series. As in the previous sections, let  $\phi_1, \dots, \phi_p$  be

the AR parameters and let  $\theta_1, \dots, \theta_q$  be the MA parameters of the  $ARMA(p, q)$  series. The polynomials used in the discussion on stationarity, invertibility, etc., in [Section 13.5](#) are

$$\theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q, \quad \phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p.$$

Let us first consider a mean zero  $MA(q)$  series which is of the form  $X_t = \varepsilon_t + \sum_{j=1}^q \theta_j \varepsilon_{t-j}$ , where  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ . So  $\{X_t\}$  is a filtered series of  $\{\varepsilon_t\}$  with  $\psi_0 = 1$ ,  $\psi_j = \theta_j, j = 1, \dots, q$ , and  $\psi_j = 0$  otherwise. Hence the frequency response function of  $\{X_t\}$  is

$$\begin{aligned} \Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi iwj) = 1 + \sum_{j=1}^q \theta_j \exp(-2\pi iwj) \\ &= \theta(z), \text{ with } z = \exp(-2\pi iw). \end{aligned}$$

Since the spectral density function of  $\{\varepsilon_t\}$  is  $f_\varepsilon(w) = \sigma^2$  for all  $w$ , the spectral density of  $\{X_t\}$  is

$$f_X(w) = |\Psi(w)|^2 f_\varepsilon(w) = \sigma^2 |\theta(z)|^2, \text{ with } z = \exp(-2\pi iw).$$

Let us now find the spectral density function of an  $AR(p)$  series. For this series

$$\begin{aligned} X_t &= \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \varepsilon_t, \text{ ie,} \\ \varepsilon_t &= X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p}. \end{aligned}$$

Thus  $\{\varepsilon_t\}$  is a filtered series of  $\{X_t\}$  with  $\psi_0 = 1$ ,  $\psi_j = -\phi_j, j = 1, \dots, p$ , and  $\psi_j = 0$  otherwise. The frequency response function is

$$\begin{aligned} \Psi(w) &= \sum_{j=-\infty}^{\infty} \psi_j \exp(-2\pi iwj) = 1 - \sum_{j=1}^p \phi_j \exp(-2\pi iwj) \\ &= \phi(z), \text{ with } z = \exp(-2\pi iw). \end{aligned}$$

It then follows that

$$\begin{aligned} f_\varepsilon(w) &= |\phi(w)|^2 f_X(w), \text{ and hence} \\ f_X(w) &= \sigma^2 \frac{1}{|\phi(z)|^2}, \text{ with } z = \exp(-2\pi iw). \end{aligned}$$

Now let us look at the  $ARMA(p, q)$  series

$$\begin{aligned} X_t &= \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}, \text{ ie,} \\ X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p} &= \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}. \end{aligned}$$

The left-hand side of the last expression is a filtered series of  $\{X_t\}$  with the frequency response function  $\Psi_1(w)$  and the right-hand side is a filtered series  $\{\varepsilon_t\}$  with frequency response function  $\Psi_2(w)$ , where  $\Psi_1(w) = \phi(z)$  and  $\Psi_2(w) = \theta(z)$  with  $z = \exp(-2\pi iw)$ . Since the spectral density function of the filtered series on the left equals the spectral density function on the right-hand side, we have

$$|\Psi_1(w)|^2 f_X(w) = |\Psi_2(w)|^2 f_\varepsilon(w) = \sigma^2 |\psi_2(w)|^2, \text{ ie,}$$

$$|\theta(z)|^2 f_X(w) = \sigma^2 |\theta(z)|^2 f_\varepsilon(w).$$

Therefore, the spectral density function of an ARMA( $p, q$ ) series is

$$f_X(w) = \sigma^2 \frac{|\theta(z)|^2}{|\theta(z)|^2}, \text{ with } z = \exp(-2\pi iw).$$

Since

$$\begin{aligned} \theta(z) &= 1 + \sum_{j=1}^q \theta_j \exp(-2\pi i j w) \\ &= 1 + \sum_{j=1}^q \theta_j \cos(2\pi j w) - i \sum_{j=1}^q \theta_j \sin(2\pi j w), \\ \phi(z) &= 1 - \sum_{j=1}^p \phi_j \exp(-2\pi i j w) \\ &= 1 - \sum_{j=1}^p \phi_j \cos(2\pi j w) + i \sum_{j=1}^p \phi_j \sin(2\pi j w), \end{aligned}$$

the spectral density function  $f_X$  of an ARMA( $p, q$ ) series can be written as

$$f_X(w) = \sigma^2 \frac{\left[1 + \sum_{j=1}^q \theta_j \cos(2\pi j w)\right]^2 + \left[\sum_{j=1}^q \theta_j \sin(2\pi j w)\right]^2}{\left[1 - \sum_{j=1}^p \phi_j \cos(2\pi j w)\right] + \left[\sum_{j=1}^p \phi_j \sin(2\pi j w)\right]^2}.$$

### Some Special Cases

#### I. Spectral density of AR(1).

When  $z = \exp(-2\pi iw)$ , we have

$$|\phi(z)|^2 = 1 + \phi^2 - 2\phi \cos(2\pi w).$$

Hence we get

$$f_X(w) = \sigma^2 \frac{1}{1 + \phi^2 - 2\phi \cos(2\pi w)}.$$

#### II. Spectral density of AR(2).

Using some algebra, we have

$$\begin{aligned} |\phi(z)|^2 &= [1 - \phi_1 \cos(2\pi w) - \phi_2 \cos(4\pi w)]^2 + [\phi_1 \sin(2\pi j w) + \phi_2 \sin(4\pi w)]^2 \\ &= 1 + \phi_1^2 + \phi_2^2 + 2\phi_1(1 - \phi_2) \cos(2\pi w) - 2\phi_2 \cos(4\pi w). \end{aligned}$$

Since  $\cos(4\pi w) = 2\cos^2(2\pi w) - 1$ , substituting this in the last expression, we have

$$|\phi(z)|^2 = \phi_1^2 + (1 + \phi_2)^2 + 2\phi_1(1 - \phi_2)\cos(2\pi w) - 4\phi_2\cos^2(2\pi w),$$

which is a quadratic polynomial in  $\cos(2\pi w)$ . Thus the spectral density function is

$$f_X(w) = \sigma^2 \frac{1}{\phi_1^2 + (1 + \phi_2)^2 + 2\phi_1(1 - \phi_2)\cos(2\pi w) - 4\phi_2\cos^2(2\pi w)}.$$

### III. Spectral density of MA(1).

Here

$$|\theta(z)|^2 = 1 + \theta^2 + 2\theta\cos(2\pi w), \text{ with } z = \exp(-2\pi iw),$$

and hence the spectral density is

$$f_X(w) = \sigma^2[1 + \theta^2 + 2\theta\cos(2\pi w)].$$

### IV. Spectral density of ARMA(1, 1).

If the parameters of an ARMA(1, 1) series are  $\phi$  and  $\theta$ , then

$$|\theta(z)|^2 = 1 + \theta^2 + 2\theta\cos(2\pi w), \text{ and}$$

$$|\phi(z)|^2 = 1 + \phi^2 - 2\phi\cos(2\pi w), \text{ with } z = \exp(-2\pi iw).$$

Thus the spectral density of an ARMA(1, 1) series is

$$f_X(w) = \sigma^2 \frac{1 + \theta^2 + 2\theta\cos(2\pi w)}{1 + \phi^2 - 2\phi\cos(2\pi w)}.$$

## Exercises

In all the problems below, the autocorrelations and autocovariances of stationary series are denoted by  $\{\gamma(h)\}$  and  $\{\rho(h)\}$ , respectively. For any ARMA series, it is understood that the mean and variance of the innovations  $\{\varepsilon_t\}$  are 0 and  $\sigma^2$ , respectively.

- 13.1.** Let  $\{X_{1t}\}, \dots, \{X_{kt}\}$  be  $k$  independent stationary series with autocovariance functions  $\gamma_1, \dots, \gamma_k$ . Then show that  $W_t = c_1X_{1t} + \dots + c_kX_{kt}$ , where  $c_1, \dots, c_k$  are constants, is also stationary. Find the autocovariances and autocorrelations of the series  $\{W_t\}$  in terms of the autocovariance and autocorrelation functions of the series  $\{X_{1t}\}, \dots, \{X_{kt}\}$ .
- 13.2.** For each of the following models, determine if it is stationary and invertible. It is understood that  $\{\varepsilon_t\}$  are iid with mean 0 and variance  $\sigma^2$ .
- (i)  $X_t = 6 + \varepsilon_t + 1.2\varepsilon_{t-1}$ .
  - (ii)  $X_t = -5 + \varepsilon_t + 0.6\varepsilon_{t-1} + 0.7\varepsilon_{t-2}$ .
  - (iii)  $X_t = 3 + 0.5X_{t-1} + \varepsilon_t = 0.4\varepsilon_{t-1}$ .
  - (iv)  $X_t = 9 + 0.7X_{t-1} + 0.6X_{t-2} + \varepsilon_t$ .
  - (v)  $X_t = 2 - 0.5X_{t-1} - 0.4X_{t-2} + \varepsilon_t + 0.3\varepsilon_{t-1} + 0.6\varepsilon_{t-2}$ .

- 13.3.** Assume that  $\{X_t\}$  follows a stationary ARMA(1, 1) model with the autoregressive and moving average parameters  $\phi$  and  $\theta$ , respectively, with  $\phi \neq -\theta$ .
- (a) Show that for any positive integer  $r$ , one may write  $X_t$  as

$$X_t = \phi^r X_{t-r} + \sum_{j=0}^{r-1} \psi_j \varepsilon_{t-j} + \phi^{r-1} \theta \varepsilon_{t-r},$$

- where  $\psi_0 = 1$  and  $\psi_j = (\phi + \theta)\phi^{j-1}$ ,  $j = 1, \dots$
- (b) Use the result in part (a) to argue that  $X_t = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}$ , where  $\{\psi_j\}$  are as given in part (a).
- (c) Show that  $\text{Var}[X_t] = [1 + (\phi + \theta)^2(1 - \phi^2)^{-1}] \sigma^2$ .
- (d) Show that the autocovariance function of  $\{X_t\}$  is
- $$\gamma(h) = (\phi + \theta)(1 + \phi\theta)(1 - \phi^2)^{-1}\sigma^2, h = 1, 2, \dots$$
- 13.4.** Let  $\{X_t\}$  be a stationary ARMA(1, 1) series with mean  $\mu$ . Assuming that  $\phi$ ,  $\theta$ , and  $\sigma^2$  are known, use the results in Exercise 13.3 to obtain an estimate of  $\tau_n^2 = \text{Var}[\hat{\mu}]$  in terms of  $\phi$ ,  $\theta$ ,  $\sigma^2$ , and  $n$ , where  $\hat{\mu} = n^{-1} \sum_{t=1}^n X_t$  is an estimate of  $\mu$  based on the available data  $X_1, \dots, X_n$ .
- 13.5.** Let  $\{X_t\}$  be stationary (not necessarily AR( $p$ )) with mean  $\mu$ . Denote  $\rho(1)$  by  $\phi$ .
- (a) Show that the best linear predictor (forecast) of  $X_t$  from  $X_{t-1}$  is
- $$X_t^{(f)} = \mu + \phi(X_{t-1} - \mu).$$
- (b) Show that the best linear predictor (backcast) of  $X_{t-2}$  from  $X_{t-1}$  is
- $$X_{t-2}^{(b)} = \mu + \phi(X_{t-1} - \mu).$$
- (c) Show that  $\text{Var}[X_t - X_t^{(f)}] = (1 + \phi^2)\gamma(0) - 2\phi\gamma(1)$ .
- (d) Show that  $\text{Var}[X_{t-2} - X_{t-2}^{(f)}] = (1 + \phi^2)\gamma(0) - 2\phi\gamma(1)$ .
- (e) Show that  $\text{Cov}[X_t - X_t^{(f)}, X_{t-2} - X_{t-2}^{(b)}] = \gamma(2) - 2\phi\gamma(1) + \phi^2\gamma(0)$ .
- (f) Show that the partial correlation between  $X_t$  and  $X_{t-2}$  given  $X_{t-1}$  is given by  $[\gamma(2) - 2\phi\gamma(1) + \phi^2\gamma(0)]/[(1 + \phi^2)\gamma(0) - 2\phi\gamma(1)]$ .
- 13.6.** Suppose that when fitting an AR( $p$ ) model to the data  $X_1, \dots, X_n$  from a stationary series  $\{X_t\}$  with mean  $\mu$ , the data are expanded by padding with  $2p$  extra values the sample mean  $\bar{X}$  at the beginning and at the end. In particular, let  $\tilde{X}_t$ ,  $t = -p + 1, \dots, n + p$  be such that  $\tilde{X}_t = X_t$  if  $1 \leq t \leq n$  and  $\tilde{X}_t = \bar{X}$  otherwise. Show that the normal equations obtained by minimizing  $\sum_{t=1}^{n+p} (Y_t - \phi_1 Y_{t-1} - \dots - \phi_p Y_{t-p})^2$ , where  $Y_t = \tilde{X}_t - \bar{X}$ , with respect to  $\phi_1, \dots, \phi_p$  are the same as the Yule-Walker equations.
- 13.7.** For an invertible MA(1) model obtain the partial correlation function  $\{\pi(h)\}$  and express it as a function of the moving average parameter  $\theta$ .
- 13.8.** Let  $\{X_t\}$  be a stationary, invertible, and nonredundant ARMA(1, 1) series with zero mean, and AR and MA parameters  $\phi$  and  $\theta$ .
- (a) Obtain an invertible representation of  $\{X_t\}$  as Eq. (6) by finding  $\{\psi_j\}$  explicitly in terms of  $\phi$  and  $\theta$ .

- (b) If the  $ARMA(1, 1)$  series  $\{X_t\}$  has mean zero and its invertible representation is  $X_{t+1} = \sum_{j=0} \pi_j X_{t-j} + \varepsilon_t$ , then consider predicting  $X_{n+1}$  using  $\hat{X}_{n+1} = \pi_1 X_n + \dots + \pi_n X_1$ . Find the mean square error of prediction  $E[\hat{X}_{n+1} - X_{n+1}]^2$ .

- 13.9.** In [Section 13.5](#) (after [Lemma 13.5.1](#)), it has been pointed out that the autocorrelation function of a stationary  $AR(2)$  series is representable in terms of the roots  $z_1, z_2$  of the equation  $g(z) = 0$ , where  $g(z) = z^2 - \phi_1 z - \phi_2$ . Condition of stationarity requires that  $(\phi_1, \phi_2)$  is inside the triangular region

$$\Delta = \{(u_1, u_2): |u_2| < 1, |u_1/(1 - u_2)| < 1\}.$$

Show that  $(\phi_1, \phi_2)$  is inside  $\Delta$  if and only if  $|z_1| < 1$  and  $|z_2| < 1$ .

- 13.10.** Let  $\{X_t\}$  be an  $AR(2)$  series with autoregressive coefficients  $\phi_1$  and  $\phi_2$ .

- (a) Show that  $|\phi_2| < 1$ .
- (b) Show that  $\rho(1) = \phi_1/(1 - \phi_2)$ .
- (c)  $\sigma^2 = (1 - \phi_1^2 - \phi_2^2)\gamma(0) - 2\phi_1\phi_2\gamma(1)$ .
- (d) Show that  $\rho(h) = \phi_1\rho(h-1) + \phi_2\rho(h-2)$ ,  $h \geq 2$ .

- 13.11.** For an  $AR(1)$  series  $\{X_t\}$  with autoregressive coefficient  $-1 < \phi < 1$ , show that the prediction error for predicting  $X_{n+h}$  by its best linear predictor based on  $X_1, \dots, X_n$  is  $\sigma^2(h) = \sigma^2 \frac{1 - \phi^{2h}}{1 - \phi^2}$ .

- 13.12.** (a) Let  $\{X_t\}$  be a mean zero time series following an  $MA(1)$  model. Let  $\hat{X}_t$  be the forecasted value of  $X_t$  based on the past  $X_{t-1}, X_{t-2}, \dots$ , and  $\hat{X}_{t+1}$  be the forecasted value of  $X_{t+1}$  based on the past  $X_t, X_{t-1}, \dots$ . Show that  $\hat{X}_{t+1} = \theta(X_t - \hat{X}_t)$ . [Here  $\theta$  is the moving average parameter.]  
(b) The series  $\{Y_t\}$  follows an  $ARIMA(0, 1, 1)$  model and assume that the series  $\{X_t\}$ , the first difference of  $\{Y_t\}$ , has zero mean. Let  $\hat{Y}_t$  be the forecasted value of  $Y_t$  based on the past  $Y_{t-1}, Y_{t-2}, \dots$ , and  $\hat{Y}_{t+1}$  be the forecasted value of  $Y_{t+1}$  based on the past  $Y_t, Y_{t-1}, \dots$ . Show that  $\hat{Y}_{t+1} = Y_t + \theta(Y_t - \hat{Y}_t)$ , where  $-1 < \theta < 1$  is the moving average coefficient.

- 13.13.** If  $\{X_t\}$  follows an  $MA(q)$  model, then show that the asymptotic variance of  $\hat{\rho}(h)$  is  $\sum_{j=-q}^q \rho(j)^2/n$ , for any  $h \geq q + 1$ .

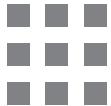
- 13.14.** (a) If  $\{X_t\}$  is stationary  $AR(1)$  and the parameter  $\phi$  is estimated using the Yule-Walker method based on a sample  $X_1, \dots, X_n$ . Obtain the asymptotic distribution of  $\sqrt{n}(\hat{\phi} - \phi)$  and explicitly obtain the parameters of this distribution in terms of  $\phi$  and  $\sigma^2$ .

- (b) Repeat part (a) for the parameter estimate  $\hat{\theta}$  of  $\theta$  for an  $MA(1)$  series.
- (c) If  $\{X_t\}$  is stationary  $ARMA(1, 1)$  and it is invertible and nonredundant. Let  $\hat{\phi}$  and  $\hat{\theta}$  be the estimates of  $\phi$  and  $\theta$  using the methods in [Section 13.8.3](#). Find the joint asymptotic distribution of  $\sqrt{n}(\hat{\phi} - \phi)$  and  $\sqrt{n}(\hat{\theta} - \theta)$ , and explicitly obtain the parameters of this distribution in terms of  $\phi, \theta$ , and  $\sigma^2$ .

- 13.15.** Prove [Lemma 13.6.1](#).

- 13.16.** Let  $\{X_t\}$  be as in Eq. (7). Show that  $\{X_t\}$  is stationary with the autocovariance function  $\{\gamma(h)\}$ , where  $\gamma(h) = \sum_{j=1}^M \sigma_j^2 \cos(2\pi w_j h)$ .

- 13.17.** Let  $\{X_t\}$  be stationary and define  $W_t = (X_{t-2} + 2X_{t-1} + 3X_t + 2X_{t+1} + X_{t+2})$ .
- (a) Find the frequency response function  $\Psi(w)$ .
  - (b) Obtain the spectral density function of  $\{W_t\}$  in terms of the spectral density function of  $\{X_t\}$ .
  - (c) Plot the square of the absolute value of the frequency response function.
- 13.18.** For a stationary series  $\{X_t\}$ , its second difference is  $W_t = X_t - 2X_{t-1} + X_{t-2}$ .
- (a) Find the frequency response function  $\Psi(w)$ .
  - (b) Obtain the spectral density function of  $\{W_t\}$  in terms of the spectral density function of  $\{X_t\}$ .
  - (c) Plot the square of the absolute value of the frequency response function.



# Appendix A

## Results From Analysis and Probability

### A.1 Some Important Results in Integration Theory

**Theorem A.1.1** (Lebesgue Dominated Convergence). *Let  $\{f_n\}$  be a sequence of integrable functions on  $\mathcal{X}$ . If*

- (i)  $\lim_{n \rightarrow \infty} f_n(x) = f(x)$  a.e. in  $\mathcal{X}$ , that is, for all  $x \notin S$  where  $\int_S dx = 0$ , and
  - (ii) there is an integrable function  $g$  on  $\mathcal{X}$  such that  $|f_n(x)| \leq g(x)$  for all  $n$  and for all  $x \in \mathcal{X}$ ,
- then

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f_n(x) dx = \int_{\mathcal{X}} \lim_{n \rightarrow \infty} f_n(x) dx = \int_{\mathcal{X}} f(x) dx.$$

*Proof.* See Royden [65, p. 88]. □

The next two theorems follow from the Dominated Convergence Theorem in a straightforward manner.

**Theorem A.1.2** (Monotone Convergence). *Let  $\{f_n\}$  be a sequence of nonnegative functions on  $\mathcal{X}$  such that  $0 \leq f_1(x) \leq f_2(x) \leq \dots$  and let  $\lim_{n \rightarrow \infty} f_n(x) = f(x)$  where  $f$  is integrable on  $\mathcal{X}$ . Then*

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f_n(x) dx = \int_{\mathcal{X}} f(x) dx.$$

**Theorem A.1.3** (Differentiation Under Integration). *Let  $f(x, t)$  for  $(x, t) \in [a, b] \times [c, d]$  be such that*

- (i)  $f(x, t)$  is an integrable function of  $x$  on  $[a, b]$  for each  $t \in [c, d]$ ,
- (ii) the partial derivative  $\partial f / \partial t$  exists and is bounded on  $[a, b] \times [c, d]$ .

Then

$$\frac{d}{dt} \int_a^b f(x, t) dx = \int_a^b \frac{\partial f(x, t)}{\partial t} dx.$$

**Theorem A.1.4** (Fubini). *If  $f(x, y)$  is integrable on  $\mathcal{X} \times \mathcal{Y}$ , then*

$$\iint_{\mathcal{X} \times \mathcal{Y}} f(x, y) dx dy = \int_{\mathcal{X}} \left[ \int_{\mathcal{Y}} f(x, y) dy \right] dx = \int_{\mathcal{Y}} \left[ \int_{\mathcal{X}} f(x, y) dx \right] dy,$$

*that is, the double integral can be evaluated interactively either way.*

## A.2 Convex Functions

**Definition A.2.1.** A real-valued function  $f$  on an interval  $(a, b) \subset \mathbb{R}$ , or more generally on  $(a_1, b_1) \times \cdots \times (a_k, b_k) \subset \mathbb{R}^k$ , is said to be convex if for any  $x_1, x_2$  in its domain and for any  $0 < \lambda < 1$ ,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2).$$

The function is strictly convex if the above inequality is a strict inequality for all such  $x_1, x_2$  and  $\lambda$ .

Geometrically, a function is convex if the straight line joining any two points on its graph lies entirely above the graph.

A twice differentiable function  $f$  on  $(a, b) \subset \mathbb{R}$  is convex iff  $f''(x) \geq 0$  for all  $a < x < b$ . More generally, a function  $f$  on  $(a_1, b_1) \times \cdots \times (a_k, b_k) \subset \mathbb{R}^k$  for which all second partial derivatives exist and are finite, is convex iff the Hessian (ie, the matrix of second partial derivatives) is nonnegative definite.

From the above definition it follows by induction, that if  $X$  is an rv taking values  $x_1, \dots, x_r$  in  $(a, b)$  with  $P[X = x_i] = \lambda_i, i = 1, \dots, r$  with  $\sum_{i=1}^r \lambda_i = 1$ , then

$$f(\mathbb{E}[X]) = f\left(\sum_{i=1}^r \lambda_i x_i\right) \leq \sum_{i=1}^r \lambda_i f(x_i) = \mathbb{E}[f(X)].$$

Obviously, this inequality also holds for a random vector  $X$  taking values in  $(a_1, b_1) \times \cdots \times (a_k, b_k) \subset \mathbb{R}^k$ .

The following theorem asserts that this inequality holds for arbitrary random vector  $X$  with finite expectation.

**Theorem A.2.1** (Jensen's Inequality). *If  $f$  is a convex function on  $I = (a_1, b_1) \times \cdots \times (a_k, b_k) \subset \mathbb{R}^k$  and  $X$  is a  $k$ -dim random vector with  $P[X \in I] = 1$ , and with finite expectation, then*

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

*Moreover the above inequality is strict if  $f$  is strictly convex, unless  $X$  is a constant with probability 1.*

*Proof.* The proof depends on the Supporting Hyperplane Theorem for convex sets. See Ferguson [1, p. 76]. □

## An Application

Consider the function  $f(x, y) = -x^p y^{1-p}$ ,  $x > 0$ ,  $y > 0$ , and  $0 < p < 1$ , which can be shown to be convex on the first quadrant on  $R^2$  by verifying that its Hessian is nonnegative definite. The following theorem is now proved by an application of Jensen's Inequality.

**Theorem A.2.2** (Hölder's Inequality). *If  $X$  and  $Y$  are positive rv's with finite means and  $0 \leq p \leq 1$ , then*

$$E[X^p Y^{1-p}] \leq (E[X])^p (E[Y])^{1-p}.$$

Taking  $X = U^2$  and  $Y = V^2$  and  $p = 1/2$ , the Cauchy-Schwarz inequality follows as a special case.

## A.3 Stieltjes Integral

Let  $f$  and  $g$  be real-valued functions on an interval  $[a, b]$ . The Stieltjes integral of  $f$  with respect to  $g$  on  $[a, b]$ , denoted by  $\int_a^b f(x) dg(x)$  or simply  $\int_a^b f dg$  is a generalization of the Riemann integral  $\int_a^b f(x) dx$ . As in Riemann integration, we need the concept of partitions.

### Definition A.3.1.

- (i) A partition of  $[a, b]$  is a finite set of real numbers  $P = \{x_0, x_1, \dots, x_n\}$  where  $a = x_0 < x_1 < \dots < x_n = b$ , of which  $[x_{i-1}, x_i]$  are segments with length  $\Delta x_i = x_i - x_{i-1}$ , and  $\Delta P = \max\{\Delta x_i, i = 1, \dots, n\}$  is the norm of  $P$ .
- (ii) A partition  $Q = \{y_0, y_1, \dots, y_m\}$  is a refinement of  $P = \{x_0, x_1, \dots, x_n\}$  if  $P \subset Q$ , in which case,  $\Delta Q \leq \Delta P$ .
- (iii) A partition  $Q = \{\xi_1, \dots, \xi_n\}$  is an intermediate partition of  $P = \{x_0, x_1, \dots, x_n\}$  if  $x_{i-1} \leq \xi_i \leq x_i$  for all  $i$ .
- (iv) For real-valued functions  $f$  and  $g$  on  $[a, b]$ , a partition  $P = \{x_0, x_1, \dots, x_n\}$  of  $[a, b]$  and an intermediate partition  $Q = \{\xi_1, \dots, \xi_n\}$  of  $P$ , the Stieltjes sum of  $f$  with respect to  $g$  on  $[a, b]$  corresponding to  $P$  and  $Q$  is defined as

$$S(f, g, P, Q) = \sum_{i=1}^n f(\xi_i) \Delta g_i, \quad \text{where } \Delta g_i = g(x_i) - g(x_{i-1}).$$

- [This generalizes the Riemann Sum  $\sum_{i=1}^n f(\xi_i) \Delta x_i$ , where  $\Delta x_i = x_i - x_{i-1}$ .]
- (v) The Stieltjes integral of  $f$  with respect to  $g$  on  $[a, b]$  is defined as a number  $\int_a^b f(x) dg(x)$  having the property that for every  $\varepsilon > 0$  there exists a  $\delta > 0$  such that

$$\left| S(f, g, P, Q) - \int_a^b f dg \right| < \varepsilon$$

for all  $P$  with  $\Delta P < \delta$  and for all intermediate partitions  $Q$  of  $P$ .

[Definition (v) is a formal way of saying,  $\int_a^b f dg = \lim_{\Delta P \rightarrow 0} S(f, g, P, Q)$ .]

As in Riemann integration,

$$\begin{aligned}\int_b^a f dg &= - \int_a^b f dg, \quad \text{for all } a < b, \text{ and} \\ \int_a^a f dg &= 0, \quad \text{by convention.}\end{aligned}$$

## Existence

The Stieltjes integral  $\int_a^b f dg$  exists if  $f$  is continuous and  $g$  is nondecreasing on  $[a, b]$ . More generally,  $\int_a^b f dg$  exists if  $f$  has at most a finite number of discontinuities,  $g$  is of bounded variation (as defined below) and  $f$  and  $g$  have no common discontinuity on  $[a, b]$ . For practical purposes, this generality will suffice.

**Definition A.3.2** (Functions of Bounded Variation). The variation of  $f$  on  $[a, b]$  corresponding to a partition  $P = \{x_0, x_1, \dots, x_n\}$

$$V^P(f) = \sum_{i=1}^n |f(x_i) - f(x_{i-1})|.$$

Obviously, if  $Q$  is a refinement of  $P$ , then  $V^P(f) \leq V^Q(f)$ . If the set  $\{V^P(f): P \text{ is a partition of } [a, b]\}$  is bounded, then  $f$  is called a function of bounded variation on  $[a, b]$  and  $V^{[a,b]}(f) = \sup_P V^P(f)$  is the total variation of  $f$  on  $[a, b]$ .

If  $f$  is differentiable on  $[a, b]$  with  $|f'(x)| \leq M$  for all  $x \in [a, b]$ , then  $f$  is of bounded variation on  $[a, b]$  and  $V^{[a,b]}(f) \leq M(b-a)$ . More generally,  $f$  is of bounded variation iff it is the difference of two nondecreasing functions. Indeed, for a function of bounded variation if we let

$$v(x; f) = V^{[a,x]}(f) \quad \text{for } a < x \leq b \text{ and } v(a; f) = 0,$$

then the function  $v(x; f)$  called the total variation function of  $f$  and the function  $r(x; f) = v(x; f) - f(x)$  called the residual function of  $f$ , on  $[a, b]$ , are both nondecreasing and  $f = v - r$ .

## Properties of Stieltjes Integrals

1. If  $f$  is Riemann integrable and  $g$  has continuous derivative  $g'$  on  $[a, b]$ , then the Riemann integral  $\int_a^b f(x) g'(x) dx$  and the Stieltjes integral  $\int_a^b f(x) dg(x)$  both exist and are equal.

- 2. (a)** If  $\int_a^b f_i \, dg$ ,  $i = 1, 2$ , exist, then so does  $\int_a^b (k_1 f_1 + k_2 f_2) \, dg$  for  $k_1$  and  $k_2$  constants and

$$\int_a^b (k_1 f_1 + k_2 f_2) \, dg = k_1 \int_a^b f_1 \, dg + k_2 \int_a^b f_2 \, dg,$$

and (a') if  $\int_a^b f d g_i$ ,  $i = 1, 2$ , exist, then so does  $\int_a^b f d(k_1 g_1 + k_2 g_2)$  for  $k_1, k_2$  constants and

$$\int_a^b f \, d(k_1 g_1 + k_2 g_2) = k_1 \int_a^b f \, d g_1 + k_2 \int_a^b f \, d g_2.$$

- (b)** If  $\int_a^b f \, dg$  exists, then for  $a < c < b$ ,  $\int_a^c f \, dg$  and  $\int_c^b f \, dg$  exist and

$$\int_a^b f \, dg = \int_a^c f \, dg + \int_c^b f \, dg.$$

- (c)** If  $\int_a^b f \, dg$  exists, then  $\int_c^d f \, dg$  exists for any  $[c, d] \subset [a, b]$ .

**Theorem A.3.1** (Integration by Parts). *If  $\int_a^b f \, dg$  exists, then  $\int_a^b g \, df$  also exists and*

$$\int_a^b f \, dg = f(b)g(b) - f(a)g(a) - \int_a^b g \, df.$$

- 4. Change of variable.** Suppose that  $\int_a^b f \, dg$  exists,  $h$  is a strictly increasing and continuous function on  $[p, q]$  with  $h(p) = a$  and  $h(q) = b$ . Then for  $F = f \circ h$  and  $G = g \circ h$  on  $[p, q]$ ,  $\int_p^q F \, dG$  exists and is equal to  $\int_a^b f \, dg$ .

*Riemann-Stieltjes (R-S) integral.* If in the Stieltjes Sum  $S(f, g, P, Q)$  we replace  $f(\xi_i)$  by  $m_i = \inf_{x \in [x_{i-1}, x_i]} f(x)$  or  $M_i = \sup_{x \in [x_{i-1}, x_i]} f(x)$ , then the resulting sums, denoted by  $\underline{RS}(f, g, P)$  and  $\overline{RS}(f, g, P)$ , respectively, are called the Lower and the Upper R-S sums which are generalizations of the Lower and Upper Riemann Sums. Since  $\underline{RS}(f, g, P) \leq \overline{RS}(f, g, Q)$  for all partitions  $P, Q$  of  $[a, b]$ , we have

$$\int_a^b f \, dg := \sup_P \underline{RS}(f, g, P) \leq \inf_P \overline{RS}(f, g, P) := \overline{\int_a^b f \, dg}.$$

If  $\int_a^b f \, dg$  and  $\overline{\int_a^b f \, dg}$ , called the lower and upper R-S integrals are equal, then the common value is called the R-S integral of  $f$  with respect to  $g$  on  $[a, b]$ .

The Stieltjes integral and the Riemann-Stieltjes integral both exist and are equal if  $f$  has at most a finite number of discontinuities,  $g$  is nondecreasing, and  $f$  and  $g$  do not have any common discontinuity on  $[a, b]$ .

*Expected value of a random variable.* The expected value of an rv with cdf  $F$  is defined as

$$E_F[X] = \int x \, dF(x).$$

More generally,  $E_F[g(X)] = \int g(x) dF(x)$  where  $g$  has at most a finite number of discontinuities.

If  $F$  is differentiable,  $F' = f$  is the pdf of  $X$ , and  $\int x dF(x) = \int xf(x) dx$ . If  $F$  increases only by jumps at  $x_1, x_2, \dots$  with jump size  $f(x_i)$  at  $x_i$ , then  $X$  is discrete with pmf  $f(x_i), i = 1, 2, \dots$ , and  $\int x dF(x) = \sum_i x_i f(x_i)$ .

If  $P[X \geq 0] = 1$ , then we have the following alternative expression for  $E_F[X] < \infty$ , using integration by parts:

$$\begin{aligned} E_F[X] &= \int_0^\infty x dF(x) = - \int_0^\infty x d[1 - F(x)] \\ &= \int_0^\infty [1 - F(x)] dx = \int_0^\infty P[X > x] dx, \end{aligned}$$

because  $x[1 - F(x)] \leq \int_x^\infty yF(y) dy$  which converges to 0 as  $x \rightarrow \infty$  by virtue of  $E_F[X] < \infty$ .

*Empirical cdf.* For a random sample  $(X_1, \dots, X_n)$  from  $F$ , the function

$$F_n(x) = n^{-1} \sum_{i=1}^n I_{(-\infty, x]}(X_i),$$

which increases by jumps of size  $1/n$  at  $X_1, \dots, X_n$  is called the empirical cdf. It follows that

$$\int g(x) dF_n(x) = n^{-1} \sum_{i=1}^n g(X_i) = \text{Sample mean of } g(X_1), \dots, g(X_n).$$

We often want to deal with  $n^{-1} \sum_{i=1}^n g(X_i) - E_F[g(X)]$ , the difference between the sample mean of  $g(X_1), \dots, g(X_n)$  and its expected value. This can be represented as

$$\begin{aligned} n^{-1} \sum_{i=1}^n g(X_i) - E_F[g(X)] &= \int g(x) dF_n(x) - \int g(x) dF(x) \\ &= \int g(x) d[F_n(x) - F(x)]. \end{aligned}$$

## A.4 Characteristic Function, Weak Law of Large Number, and Central Limit Theorem

**Theorem A.4.1** (Helly-Bray Theorem).  $X_n \xrightarrow{\mathcal{L}} X$  implies  $E[g(X_n)] \rightarrow E[g(X)]$  for all bounded and a.e. continuous functions  $g$ .

*Proof.* See Breiman [42, p. 160]. □

**Note.** Let  $D_g$  denote the set of discontinuity points of  $g$ . If  $P[X \in D_g] = 0$ , then  $g$  is a.e. continuous.

## Characteristics Function

**Definition A.4.1.** The characteristic function (cf) of an rv  $X$  with cdf  $F$ , or the cf of  $F$  is defined to be

$$\varphi(t) = \varphi_X(t) = E[e^{itX}] = \int_{-\infty}^{\infty} e^{itx} dF(x) = \int_{-\infty}^{\infty} \cos(tx) dF(x) + i \int \sin(tx) dF(x),$$

where  $t$  is real and  $i = \sqrt{-1}$ . In general,  $\varphi(t)$  is complex-valued, with

$$\varphi(0) = 1, |\varphi(t)| \leq \int_{-\infty}^{\infty} dF(x) = 1 \text{ for all } t, \text{ and } \varphi(-t) = \overline{\varphi(t)} \text{ for all } t,$$

where  $\bar{z}$  denotes the complex conjugate of a complex number  $z$ .

If  $f$  is symmetrically distributed (about 0) (ie,  $X \stackrel{\mathcal{D}}{=} -X$ ), then

$$\overline{\varphi(t)} = \varphi(-t) = E[e^{i(-t)X}] = E[e^{it(-X)}] = E[e^{itX}] = \varphi(t).$$

Thus the cf of a symmetric rv is real-valued function of  $t$ .

By dominated convergence,  $\varphi(t)$  is continuous and if  $m_k = E[X^k]$  exists, then  $\varphi(t)$  is  $k$ -times differentiable; moreover, we can differentiate under the integral sign, that is,

$$\begin{aligned}\varphi^{(r)}(t) &= i^r \int_{-\infty}^{\infty} x^r e^{itx} dF(x) \text{ for } 0 \leq r \leq k, \text{ and} \\ \varphi^{(r)}(0) &= i^r \int_{-\infty}^{\infty} x^r dF(x) = i^r E[X^r] = i^r m_r, \quad 0 \leq r \leq k.\end{aligned}$$

In the neighborhood of  $t = 0$ , we have the McLaurin Series

$$\varphi(t) = 1 + \sum_{r=1}^k (m_r/r!) (it)^r + o(t^k), \quad \text{as } t \rightarrow 0.$$

## Special Cases

- Let  $\Phi$  denote the cdf of  $N(0, 1)$ . Then the cf of  $\Phi$  is

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} d\Phi(x) = e^{-t^2/2}.$$

- Let  $F$  denote the cdf of an rv  $X$  with  $P[X = c] = 1$  (ie,  $F(x) = 0$  for  $x < c$  and  $F(x) = 1$  for  $x \geq c$ ). Then the cf of  $F$  is  $\varphi(t) = e^{itc}$ .

## Properties of cf

- If  $\varphi_X$  is the cf of  $X$ , then  $\varphi_{aX+b}(t) = e^{itb} \varphi_X(at)$  for constant  $a, b$ .
- If  $X_1, \dots, X_n$  are independent, then  $\varphi_{\sum_{i=1}^n X_i}(t) = \prod_{i=1}^n \varphi_{X_i}(t)$ .

In particular, if  $X_1, \dots, X_n$  are iid as  $X$ , then

$$\varphi_{\sum_{i=1}^n X_i / \sqrt{n}}(t) = \{\varphi_X(t/\sqrt{n})\}^n.$$

We now state the following two fundamental theorems in the theory of characteristic functions, for the proof of which we refer to Cramér [18, p. 93–8].

**Theorem A.4.2** (Uniqueness of Characteristic Functions on  $\mathbb{R}$ ). *If  $(a - h, a + h)$  is a continuity interval of a cdf  $F$  (ie,  $a \pm h$  are continuity points of  $F$ ), and if  $\varphi$  is the cf of  $F$ , then*

$$F(a+h) - F(a-h) = \lim_{T \rightarrow \infty} \frac{1}{\pi} \int_{-T}^T \frac{\sin(ht)}{t} e^{-ita} \varphi(t) dt.$$

Consequently, the cf  $\varphi$  determines the cdf  $F$ .

**Theorem A.4.3** (Continuity of Characteristic Functions on  $\mathbb{R}$ ). *If  $\{F_n\}$  is a sequence of cdf's and  $\{\varphi_n\}$  is the corresponding sequence of cf's, then  $F_n \rightarrow F$  (at all continuity points of  $F$ ) iff there exists a  $\varphi$  which is continuous at  $t = 0$ , such that  $\varphi_n(t) \rightarrow \varphi(t)$  for all  $t$ . Moreover, if there is such a  $\varphi$ , then it is the cf of  $F$ .*

**Note.** The Uniqueness and Continuity Theorems for Characteristic Functions also extend to  $\mathbb{R}^k$ . See Cramér [18, p. 100–3].

*Remark A.4.1.* Combining the Helly-Bray Theorem with the Continuity Theorem for Characteristic Functions, we now conclude that the following are equivalent:

1.  $X_n \xrightarrow{\mathcal{L}} X$ .
2.  $F_n(x) \rightarrow F(x)$  at all continuity points of  $F$ .
3.  $E[g(X_n)] \rightarrow E[g(X)]$  for all bounded and a.e. continuous functions  $g$ .
4.  $\varphi_n(t) = E[e^{itX_n}] \rightarrow \varphi(t) = E[e^{itX}]$  for all  $t$ .

## Applications of the Continuity Theorem for Characteristic Functions

**Theorem A.4.4.** *If  $X_n \xrightarrow{\mathcal{L}} X$  and  $g$  is a a.e. continuous function, then  $g(X_n) \xrightarrow{\mathcal{L}} g(X)$ .*

*Proof.* Since  $\cos(tg(x))$  and  $\sin(tg(x))$  are a.e. continuous functions of  $x$  for every  $t$ , it follows from the Helly-Bray Theorem that

$$\begin{aligned} \varphi_{g(X_n)}(t) &= E[e^{itg(X_n)}] = E[\cos(tg(X_n))] + iE[\sin(tg(X_n))] \\ &\rightarrow E[\cos(tg(X))] + iE[\sin(tg(X))] = E[e^{itg(X)}] = \varphi_{g(X)}(t). \end{aligned}$$

Hence  $g(X_n) \xrightarrow{\mathcal{L}} g(X)$  by the Continuity Theorem. □

We now prove [Theorems 3.2.1–3.2.3 of Chapter 3](#).

**Theorem A.4.5** (Weak Law of Large Numbers (Khinchine)). *If  $X_1, X_2, \dots$  are iid as  $X$  with  $E[X] = \mu$ , then  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i \xrightarrow{\mathcal{P}} \mu$ .*

*Proof.* The cf of  $\bar{X}_n$  is

$$\varphi_{\bar{X}_n}(t) = \varphi_{n^{-1} \sum_{i=1}^n X_i}(t) = \{\varphi_X(t/n)\}^n = [1 + i(t/n)\mu + o(1/n)]^n$$

by McLaurin's series expansion of  $\varphi_X(t/n)$  about 0. Hence

$$\lim_{n \rightarrow \infty} \varphi_{\bar{X}_n}(t) = \lim_{n \rightarrow \infty} [1 + i(t/n)\mu + o(1/n)]^n = e^{it\mu},$$

which is the cf of an rv  $X_0$  which takes the value  $\mu$  with probability 1. Hence  $\bar{X}_n \xrightarrow{\mathcal{L}} X_0$  by the Continuity Theorem, so that for all  $\varepsilon > 0$ ,

$$\begin{aligned} P[\bar{X}_n \leq \mu - \varepsilon] &\rightarrow P[X_0 \leq \mu - \varepsilon] = 0 \text{ and} \\ P[\bar{X}_n \leq \mu + \varepsilon] &\rightarrow P[X_0 \leq \mu + \varepsilon] = 1. \end{aligned}$$

Thus, for any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P[|\bar{X}_n - \mu| \leq \varepsilon] = \lim_{n \rightarrow \infty} P[\bar{X}_n \leq \mu + \varepsilon] - \lim_{n \rightarrow \infty} P[\bar{X}_n \leq \mu - \varepsilon] = 1 = 0 = 1.$$

□

**Theorem A.4.6** (Central Limit Theorem (Lindeberg-Lévy)). *If  $X_1, X_2, \dots$  are iid as  $X$  with  $E[X] = \mu$  and  $\text{Var}[X] = \sigma^2$ , then*

$$Z_n = \frac{1}{\sigma\sqrt{n}} \sum_{j=1}^n (X_j - \mu) \xrightarrow{\mathcal{L}} Z \sim N(0, 1), \text{ ie, } \frac{1}{\sqrt{n}} \sum_{j=1}^n (X_j - \mu) \xrightarrow{\mathcal{L}} \sigma Z \sim N(0, \sigma^2).$$

*Proof.* Let  $(X_j - \mu)/\sigma = Y_j$ . Then  $Y_1, Y_2, \dots$  are iid as  $Y = (X - \mu)/\sigma$  with  $E[Y] = 0$  and  $E[Y^2] = 1$ . Then

$$\begin{aligned} \varphi_{Z_n}(t) &= \varphi_{n^{-1/2} \sum_{j=1}^n Y_j}(t) = \{\varphi_Y(t/\sqrt{n})\}^n \\ &= [1 + i(t/\sqrt{n})0 + (1/2)t^2(t/\sqrt{n})^2 1 + o(1/n)]^n \\ &= [1 - t^2/(2n) + o(1/n)]^n \rightarrow e^{-t^2/2} = \varphi_Z(t). \end{aligned}$$

Hence  $Z_n \xrightarrow{\mathcal{L}} Z$  by the Continuity Theorem.

□

**Theorem A.4.7** (The Cramér-Wold Device). *Let  $\{X_n\}$  be  $k$ -dim random vectors such that  $\mathbf{a}^T X_n \xrightarrow{\mathcal{L}} \mathbf{a}^T \mathbf{X}$  for all  $\mathbf{a} \in \mathbb{R}^k$ . Then  $X_n \xrightarrow{\mathcal{L}} \mathbf{X}$ .*

*Proof.* Since  $\mathbf{a}^T X_n \xrightarrow{\mathcal{L}} \mathbf{a}^T \mathbf{X}$  for all  $\mathbf{a} \in \mathbb{R}^k$ ,

$$\varphi_{\mathbf{a}^T X_n}(t) = E[e^{it(\mathbf{a}^T X_n)}] \rightarrow \varphi_{\mathbf{a}^T \mathbf{X}}(t) = E[e^{it(\mathbf{a}^T \mathbf{X})}], \quad \text{for all } t \in R \text{ and } \mathbf{a} \in \mathbb{R}^k.$$

But  $\{t\mathbf{a}: t \in \mathbb{R}, \mathbf{a} \in \mathbb{R}^k\} = \{\mathbf{s}: \mathbf{s} \in \mathbb{R}^k\}$ , and therefore,

$$\varphi_{X_n}(\mathbf{s}) = E[e^{i\mathbf{s}^T X_n}] \rightarrow E[e^{i\mathbf{s}^T \mathbf{X}}] = \varphi_{\mathbf{X}}(\mathbf{s}) \quad \text{for all } \mathbf{s} \in \mathbb{R}^k.$$

Hence  $X_n \xrightarrow{\mathcal{L}} \mathbf{X}$  by the Continuity Theorem for Characteristic Functions on  $\mathbb{R}^k$ .

□

**Theorem A.4.8** (Multivariate Central Limit Theorem). *If  $\mathbf{X}_1, \mathbf{X}_2, \dots$  are iid as  $\mathbf{X}$  in  $\mathbb{R}^k$  with the mean vector  $E[\mathbf{X}] = \boldsymbol{\mu}$  and the covariance matrix  $E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \boldsymbol{\Sigma}$  which is positive definite, then*

$$n^{-1/2} \sum_{j=1}^n (\mathbf{X}_j - \boldsymbol{\mu}) \xrightarrow{\mathcal{L}} \mathbf{W} \sim N_k(\mathbf{0}, \boldsymbol{\Sigma}).$$

*Proof.* For  $\mathbf{a} \in \mathbb{R}^k$ , let  $Y_j = \mathbf{a}^T \mathbf{X}_j$  and  $Y = \mathbf{a}^T \mathbf{X}$ . Then  $Y_1, Y_2, \dots$  are iid as  $Y$  where  $E[Y] = \mathbf{a}^T \boldsymbol{\mu}$  and  $\text{Var}[Y] = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$ , so that

$$n^{-1/2} \sum_{j=1}^n (\mathbf{a}^T \mathbf{X}_j - \mathbf{a}^T \boldsymbol{\mu}) \xrightarrow{\mathcal{L}} Z_{\mathbf{a}} \sim N(0, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a})$$

by the univariate Central Limit Theorem. On the other hand, if  $\mathbf{W} \sim N_k(\mathbf{0}, \boldsymbol{\Sigma})$ , then  $\mathbf{a}^T \mathbf{W} \sim N(0, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a})$ . Thus

$$\mathbf{a}^T \left\{ n^{-1/2} \sum_{j=1}^n (\mathbf{X}_j - \boldsymbol{\mu}) \right\} \xrightarrow{\mathcal{L}} \mathbf{a}^T \mathbf{W} \quad \text{for all } \mathbf{a} \in \mathbb{R}^k,$$

and so the theorem follows by the Cramér-Wold device.  $\square$

## A.5 Weak Convergence of Probabilities on $C[0,1]$

A metric on a nonempty set  $S$  is a function  $d: S \times S \rightarrow \mathbb{R}$  such that (i)  $d(x, y) \geq 0$  and  $d(x, y) = 0$  iff  $x = y$ , (ii)  $d(x, y) = d(y, x)$ , (iii)  $d(x, y) \leq d(x, z) + d(y, z)$  for all  $x, y, z \in S$ . We call  $(S, d)$  a metric space.

For each  $x \in S$  and  $\varepsilon > 0$ , the set  $S_{\varepsilon}(x) = \{y: d(x, y) < \varepsilon\}$  is the open sphere of radius  $\varepsilon$  centered at  $x$ . A set  $G \subset S$  is open  $\iff$  for each  $x \in G$ ,  $S_{\varepsilon}(x) \subset G$  for some  $\varepsilon > 0 \iff G$  is a union of open spheres.

In  $S$ , a sequence  $\{x_n\}$  converges to  $x$  if  $\lim_{n \rightarrow \infty} d(x_n, x) = 0$ . The boundary of  $A \subset S$  is  $\partial A = \{x \in S: x \text{ is a limit point of sequences in both } A \text{ and } A^c\}$ .

A collection of open sets  $\{G_{\lambda}: \lambda \in \Lambda\}$  is an open covering of  $A \subset S$  if  $A \subset \bigcup_{\lambda \in \Lambda} G_{\lambda}$ . A set  $K \subset S$  is compact if for every open covering of  $K$ , there is a finite subcovering.

Let  $\{P_n\}$  and  $\{P\}$  be probabilities on  $(S, \mathcal{S})$  where  $\mathcal{S}$  is the Borel  $\sigma$ -field, that is, the smallest  $\sigma$ -field of subsets of  $S$  which includes all open sets and let  $C(S)$  be the set of all bounded continuous functions  $f: S \rightarrow \mathbb{R}$ .

The sequence  $\{P_n\}$  converges weakly to  $P$  iff  $\int f \, dP_n \rightarrow \int f \, dP$  for all  $f \in C(S)$ . This is denoted by  $P_n \xrightarrow{w} P$ . Equivalently,  $P_n \xrightarrow{w} P$  iff  $P_n(A) \rightarrow P(A)$  for all  $A \in \mathcal{S}$  for which  $P(\partial A) = 0$ . The weak limit is unique (ie, if  $P_n \xrightarrow{w} P$  and  $P_n \xrightarrow{w} Q$ , then  $P = Q$ ).

A transformation  $g: (S, \mathcal{S}) \rightarrow (S', \mathcal{S}')$  is measurable if  $g^{-1}(B) \in \mathcal{S}$  for all  $B \in \mathcal{S}'$ , in which case, for  $\{P_n\}$  and  $P$  on  $(S, \mathcal{S})$ ,  $\{P_n g^{-1}\}$  and  $P g^{-1}$  are induced probabilities on  $(S', \mathcal{S}')$ , where

$$P_n g^{-1}(B) = P_n(g^{-1}(B)) \text{ and } P g^{-1}(B) = P(g^{-1}(B)) \quad \text{for all } B \in \mathcal{S}'.$$

A function  $g: (S, d) \rightarrow (S', d')$  is continuous iff  $x_n \rightarrow x$  in  $S \implies g(x_n) \rightarrow g(x)$  in  $S'$ .

**Theorem A.5.1** (Continuous Mapping Theorem). *If  $P_n \xrightarrow{w} P$  on  $(S, \mathcal{S})$  and if  $g: (S, d) \rightarrow (S', d')$  is continuous, then  $P_ng^{-1} \xrightarrow{w} Pg^{-1}$  on  $(S', \mathcal{S}')$ . More generally, the theorem holds if  $P(D_g) = 0$  where  $D_g$  is the set of discontinuity points of  $g$ .*

The continuous mapping theorem stated above is a generalization of [Theorem 3.2.5\(III\)](#) dealing with probability distributions of  $k$ -dim rv's.

Let  $C = C[0, 1]$  be the set of all continuous functions on  $[0, 1]$  with the metric  $d(x, y) = \sup_{0 \leq t \leq 1} |x(t) - y(t)|$  and let  $\mathcal{C}$  denote the Borel  $\sigma$ -field in  $C$ . On  $(C, \mathcal{C})$ , we now consider the weak convergence of  $\{P_n\}$  to  $P$ , which holds under two conditions, namely, convergence of finite-dimensional distributions (fdd) and “tightness” of the sequence  $\{P_n\}$ . For notational simplicity, we describe these conditions in terms of random elements  $X_n(\cdot)$  of  $(C, \mathcal{C}, P_n)$  and  $X(\cdot)$  of  $(C, \mathcal{C}, P)$ .

**Definition A.5.1.** For each positive integer  $k$  and  $t_1 < \dots < t_k$  in  $[0, 1]$ , the distributions of  $(X_n(t_1), \dots, X_n(t_k))$  and  $(X(t_1), \dots, X(t_k))$  are called fdd of  $X_n(\cdot)$  and  $X(\cdot)$ , respectively.

**Definition A.5.2.** A family of probabilities  $\Pi$  on  $(S, \mathcal{S})$  is tight if for every  $\varepsilon > 0$ , there exists a compact set  $K$  such that  $P(K) > 1 - \varepsilon$  for all  $P \in \Pi$ . In particular, a sequence of probabilities  $\{P_n\}$  on  $(C, \mathcal{C})$  is tight if for every  $\varepsilon > 0$ , there exists a compact set  $K$  in  $C$  such that  $P[X_n(\cdot) \in K] > 1 - \varepsilon$  for all  $n$ .

## Notation

We are writing  $P[X_n(\cdot) \in K]$  for  $P_n(K)$ .

**Theorem A.5.2.** Let  $\{X_n(\cdot), n = 1, 2, \dots\}$  and  $X(\cdot)$  denote random elements of  $(C, \mathcal{C}, P_n)$  and  $(C, \mathcal{C}, P)$ , respectively. Then  $P_n \xrightarrow{w} P$  or equivalently,  $X_n \xrightarrow{w} X$  if

- (i) the fdd's of  $\{X_n(\cdot)\}$  converge to those of  $X(\cdot)$ ,
- (ii-a)  $\{X_n(0)\}$  is tight (ie,  $X_n(0) = O_P(1)$ ), and
- (ii-b) there exist constants  $\gamma \geq 0$  and  $\alpha > 1$  and a nondecreasing, continuous function  $F$  on  $[0, 1]$  such that

$$P[|X_n(t_2) - X_n(t_1)| \geq \lambda] \leq \frac{1}{\lambda^\gamma} |F(t_2) - F(t_1)|^\alpha$$

holds for all  $t_1, t_2$  and  $n$ , and all  $\lambda > 0$  (see [43, p. 95–6]).



# Appendix B

## Basic Results From Matrix Algebra

This appendix lists some basic definitions, formulas, and results for vectors and matrices which are used in this book. We begin with some simple definitions and elementary results.

### B.1 Some Elementary Facts

It is known from the theory of matrices that the number of linearly independent rows of a matrix  $A$  equals the number of linearly independent columns, and the rank of  $A$  (denoted by  $\text{rank}(A)$ ) is defined to be the number of linearly independent rows of  $A$  (or the number of linearly independent columns). For any vector  $x$ ,  $x^T x$  will be denoted by  $\|x\|^2$ , which equals the square of the length of  $x$ . A matrix  $A$  of order  $n \times m$  is said to have a full rank if  $\text{rank}(A) = \min(n, m)$ .

For any  $n \times n$  matrix  $A$ , its quadratic form is defined to be  $q(x) = x^T A x$ , where  $x \in \mathbb{R}^n$ . If  $A$  is not symmetric, then  $q(x)$  may also be written as  $x^T \tilde{A} x$ , where  $\tilde{A} = (1/2)(A + A^T)$  is the symmetrized version of  $A$ .

**Definition B.1.1.** All the matrices in this definition are assumed to be square of order  $n$  (ie, the matrices have  $n$  rows and  $n$  columns).

- (a) Trace of a matrix is defined to be the sum of its diagonal elements (ie,  $\text{trace}(A) = \sum_{i=1}^n a_{ii}$ ).
- (b) The determinant  $A$  (denoted by  $|A|$ ) is defined to be  $\sum_{\pi} (-1)^{\pi} a_{i,\pi(i)}$ , where the sum is over all permutations  $\pi$  of  $\{1, \dots, n\}$ , and  $(-1)^{\pi}$  equals 1 or  $-1$  depending on whether  $\pi$  is a positive or a negative permutation.
- (c) A symmetric matrix  $A$  is called nonnegative definite if its quadratic form  $x^T A x \geq 0$  for any  $x \in \mathbb{R}^n$ . If  $x^T A x > 0$  for all  $0 \neq x \in \mathbb{R}^n$ , then  $A$  is called a positive definite matrix.
- (d) A matrix  $A$  is said to be orthogonal if its rows are orthonormal (ie, the row vectors are orthogonal to each other and each has unit length). Consequently,  $A A^T = I$ . It is easy to check that  $A$  is nonsingular,  $A^T = A^{-1}$  and  $A^T A = I$ . Since  $A^T A = I$ , columns of  $A$  are orthonormal.

Here are some important results on rank and trace of matrices.

**Lemma B.1.1.**

- (a) For any matrix  $A$  of order  $n \times m$ ,  $\text{rank}(AA^T) = \text{rank}(A)$ .
- (b) If  $A$  and  $B$  are of order  $n \times m$ , then  $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$ .
- (c) If  $A$  and  $B$  are of order  $n \times m$  and  $m \times k$ , respectively, then  
 $\text{rank}(AB) \leq \min[\text{rank}(A), \text{rank}(B)]$ .
- (d) If  $A$  and  $B$  are of order  $n \times n$ , then  $\text{trace}(A + B) = \text{trace}(A) + \text{trace}(B)$ .
- (e) If  $A$  and  $B$  are of order  $n \times m$  and  $m \times n$ , then  $\text{trace}(AB) = \text{trace}(BA)$ .
- (f) If  $A$  and  $B$  are of order  $n \times n$ , then  $|AB| = |A| |B|$ .

A class of formulas known as the *Sherman-Morrison Formulas* are quite useful in inverting matrices.

**Theorem B.1.1.**

- (a) If  $a \in \mathbb{R}^n$ , then  $(I + aa^T)^{-1} = I - (1 + \|a\|^2)^{-1}aa^T$ .
- (b) If  $a \in \mathbb{R}^n$ ,  $\|a\| \neq 1$ , then  $(I - aa^T)^{-1} = I + (1 - \|a\|^2)^{-1}aa^T$ .
- (c) If  $A$  is of order  $n \times m$ , then  $(I + AA^T)^{-1} = I - A(I + A^TA)^{-1}A^T$ .
- (d) If  $A$  is  $n \times m$  and  $B$  is a positive definite matrix of order  $n \times n$ , then  
 $(B + AA^T)^{-1} = B^{-1} - B^{-1}A(I + A^TB^{-1}A)^{-1}A^TB^{-1}$ .

The following is the Cauchy-Schwarz inequality for the matrices.

**Theorem B.1.2.** Let  $a$  be an  $n$ -dim vector and  $A$  be a positive definite matrix of order  $n \times n$ .

- (a) For any  $x \in \mathbb{R}^n$ ,  $|a^T x|^2 \leq \|a\|^2 \|x\|^2$ . Moreover,  $\sup\{|a^T x|^2 / \|x\|^2 : x \in \mathbb{R}^n\} = \|a\|^2$  and this supremum is attained at  $x = a$ .
- (b) For any  $x \in \mathbb{R}^n$ ,  $|a^T x|^2 \leq [a^T A^{-1} a][x^T A x]$ . Moreover,  
 $\sup\{|a^T x|^2 / [x^T A x] : x \in \mathbb{R}^n\} = a^T A^{-1} a$  and this supremum is attained at  $x = A^{-1} a$ .

## B.2 Eigenvalues and Eigenvectors

For a square matrix  $A$  of order  $n$ , if there exists a scalar (may be complex) and a vector  $x$  (may be complex) such that  $Ax = \lambda x$ , then  $\lambda$  is called an eigenvalue of  $A$  with the corresponding eigenvector  $x$ . The following result lists some basic properties of eigenvalues.

**Lemma B.2.1.**

- (a) If  $A$  is symmetric, then all its eigenvalues and eigenvectors are real.
- (b) If  $A$  is nonsingular, then all its eigenvalues are nonzero.
- (c) If  $\lambda$  is an eigenvalue of symmetric matrix  $A$  with the corresponding eigenvector  $x$ , then for any positive integer  $r$ ,  $\lambda^r$  is an eigenvalue of  $A^r$  with eigenvector  $x$ .
- (d) The nonzero eigenvalues of  $A^T A$ , where  $A$  is a matrix of order  $n \times m$ , are the same as those of  $AA^T$ .
- (e) The eigenvalues of a nonnegative definite (positive definite) matrix  $A$  are nonnegative (positive).

The following is an important result that is widely used in Linear Models and Multivariate Analysis.

**Theorem B.2.1** (Spectral Decomposition Theorem). *If  $\mathbf{A}$  is a symmetric matrix of order  $n \times n$ , then there exist an  $n \times n$  orthogonal matrix  $\mathbf{U}$  with columns  $\mathbf{u}_1, \dots, \mathbf{u}_n$  (ie,  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ ) and a diagonal matrix  $\Lambda$  of order  $n \times n$  with diagonal elements  $\lambda_1, \dots, \lambda_n$  such that*

$$\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^T = \sum_{j=1}^n \lambda_j \mathbf{u}_j \mathbf{u}_j^T.$$

Here  $\{\lambda_j\}$  are the eigenvalues of  $\mathbf{A}$  with the corresponding orthonormal eigenvectors  $\{\mathbf{u}_j\}$ .

Is there an analog of the Spectral Decomposition Theorem for an arbitrary matrix  $\mathbf{A}$  is of order  $n \times m$ ? The answer is yes. Positive square roots of the eigenvalues of  $\mathbf{A}^T\mathbf{A}$  are called the *singular values* of  $\mathbf{A}$ . Since  $\mathbf{A}^T\mathbf{A}$  is nonnegative definite, its eigenvalues are nonnegative and thus the square roots of the eigenvalues of  $\mathbf{A}^T\mathbf{A}$  are real. Clearly, if  $\mathbf{A}$  is symmetric, then  $\mathbf{A}^T\mathbf{A} = \mathbf{A}^2$  with eigenvalues  $\{\lambda_j^2\}$ , where  $\{\lambda_j\}$  are the eigenvalues of  $\mathbf{A}$ , and the singular values of  $\mathbf{A}$  are  $\{|\lambda_j|\}$ .

**Theorem B.2.2** (Singular Value Decomposition). *Let  $\mathbf{A}$  be a matrix of order  $n \times m$ ,  $m \leq n$ . There exist an  $n \times m$  suborthogonal matrix  $\mathbf{U}$  (ie, the columns of  $\mathbf{U}$  are orthonormal), an  $m \times m$  orthogonal matrix  $\mathbf{V}$ , and an  $m \times m$  diagonal matrix  $\Lambda$  with nonnegative diagonal entries  $\lambda_1, \dots, \lambda_m$  such that*

$$\mathbf{A} = \mathbf{U}\Lambda\mathbf{V}^T = \sum_{j=1}^m \lambda_j \mathbf{u}_j \mathbf{v}_j^T,$$

where  $\{\mathbf{u}_j\}$  and  $\{\mathbf{v}_j\}$  are the columns of  $\mathbf{U}$  and  $\mathbf{V}$ , respectively.

It is important to note that  $\{\lambda_j\}$  in the Singular Value Decomposition of  $\mathbf{A}$  are the singular values of  $\mathbf{A}$ , and  $\{\lambda_j^2\}$  are the eigenvalues of the matrix  $\mathbf{A}^T\mathbf{A}$ .

We finally write down two important formulas for the trace and determinant of matrices.

**Lemma B.2.2.** *Let  $\mathbf{A}$  be a symmetric matrix of order  $n \times n$  with eigenvalues  $\{\lambda_j, j = 1, \dots, n\}$ . Then,*

$$\text{trace}(\mathbf{A}) = \sum_{j=1}^n \lambda_j \quad \text{and} \quad |\mathbf{A}| = \prod_{j=1}^n \lambda_j.$$

The following result on the optimization of quadratic forms involves the eigenvalues and eigenvectors.

**Theorem B.2.3.** *Let  $\lambda_1 \geq \dots \geq \lambda_n$  be the eigenvalues of a symmetric matrix  $\mathbf{A}$  of order  $n \times n$  with the corresponding orthonormal eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_n$ . It is understood that  $\mathbf{x}, \mathbf{x}_1, \dots$  written below are in  $\mathbb{R}^n$ . Then the following hold:*

- (a)  $\sup\{\mathbf{x}^T \mathbf{A} \mathbf{x} : \|\mathbf{x}\| = 1\} = \sup_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \lambda_1$ , and this supremum is attained at  $\mathbf{x} = \mathbf{u}_1$ .
- (b)  $\inf\{\mathbf{x}^T \mathbf{A} \mathbf{x} : \|\mathbf{x}\| = 1\} = \inf_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \lambda_n$ , and this infimum is attained at  $\mathbf{x} = \mathbf{u}_n$ .

(c) For  $1 < m \leq n$ , we have

$$\begin{aligned} & \sup \left\{ \mathbf{x}_1^T \mathbf{A} \mathbf{x}_1 + \cdots + \mathbf{x}_m^T \mathbf{A} \mathbf{x}_m : \mathbf{x}_1, \dots, \mathbf{x}_m \text{ are orthonormal} \right\} \\ &= \sup \left\{ \frac{\mathbf{x}_1^T \mathbf{A} \mathbf{x}_1}{\mathbf{x}_1^T \mathbf{x}_1} + \cdots + \frac{\mathbf{x}_m^T \mathbf{A} \mathbf{x}_m}{\mathbf{x}_m^T \mathbf{x}_m} : \mathbf{x}_1 \neq \mathbf{0}, \dots, \mathbf{x}_m \neq \mathbf{0} \text{ are orthogonal} \right\} \\ &= \lambda_1 + \cdots + \lambda_m, \end{aligned}$$

and this supremum is attained at  $\mathbf{x}_1 = \mathbf{u}_1, \dots, \mathbf{x}_m = \mathbf{u}_m$ .

### B.3 Functions of Symmetric Matrices

If  $\mathbf{A}$  is a symmetric matrix of order  $n \times n$ , then for any real-valued function  $f$  whose domain contains the eigenvalues of  $\mathbf{A}$ , it is possible to define the corresponding function of  $\mathbf{A}$  using the Spectral Decomposition Theorem. More formally, let  $f$  be a real-valued function with domain  $D \subset \mathbb{R}$ , and let  $\{\lambda_j : j = 1, \dots, n\}$  be the eigenvalues of  $\mathbf{A}$  with the corresponding orthonormal eigenvector  $\{\mathbf{u}_j\}$ . If the eigenvalues of  $\mathbf{A}$  are inside the set  $D$ , then the matrix function  $f(\mathbf{A})$  is defined to be  $f(\mathbf{A}) = \sum_{j=1}^n f(\lambda_j) \mathbf{u}_j \mathbf{u}_j^T$ . Here are some examples that are useful in Linear Models and Multivariate Analysis.

- I. (Square root of a matrix) Let  $\mathbf{A}$  be nonnegative definite and let  $f(u) = u^{1/2}$ ,  $u \geq 0$  (ie,  $D = [0, \infty)$ ). Then

$$\mathbf{A}^{1/2} = f(\mathbf{A}) = \sum_{j=1}^n f(\lambda_j) \mathbf{u}_j \mathbf{u}_j^T = \sum_{j=1}^n \lambda_j^{1/2} \mathbf{u}_j \mathbf{u}_j^T.$$

Clearly,  $\mathbf{A}^{1/2}$  is symmetric and it is fairly easy to check that  $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$ .

- II. (Inverse of a matrix) If  $\mathbf{A}$  is nonsingular, then all its eigenvalues are nonzero and let  $f(u) = u^{-1}$ ,  $u \neq 0$  (ie,  $D = \mathbb{R} - \{0\}$ ). Then

$$\mathbf{A}^{-1} = f(\mathbf{A}) = \sum_{j=1}^n \lambda_j^{-1} \mathbf{u}_j \mathbf{u}_j^T.$$

It is fairly easy to check that  $\mathbf{A}^{-1}$  is indeed the inverse of the matrix  $\mathbf{A}$ .

- III. (Inverse of square root of a matrix) Let  $\mathbf{A}$  be positive definite and let  $f(u) = u^{-1/2}$ ,  $u > 0$  (ie,  $D = (0, \infty)$ ). Then

$$\mathbf{A}^{-1/2} = f(\mathbf{A}) = \sum_{j=1}^n \lambda_j^{-1/2} \mathbf{u}_j \mathbf{u}_j^T.$$

It is clear that  $\mathbf{A}^{-1/2}$  is symmetric, and it is easy to verify that  $\mathbf{A}^{-1/2} \mathbf{A}^{-1/2} = \mathbf{A}^{-1}$ .

## B.4 Generalized Eigenvalues

Let  $\mathbf{A}$  and  $\mathbf{B}$  be symmetric matrices of order  $n \times n$  where  $\mathbf{B}$  is positive definite. We say that  $\lambda$  is a (generalized) eigenvalue of  $\mathbf{A}$  with respect to  $\mathbf{B}$  if there is a vector  $\mathbf{l}$  in  $\mathbb{R}^n$  such that  $\mathbf{Al} = \lambda \mathbf{Bl}$ . Premultiplying both sides by  $\mathbf{B}^{-1}$  we get  $\mathbf{B}^{-1}\mathbf{Al} = \lambda \mathbf{l}$ . In other words, if  $\lambda$  is an eigenvalue of  $\mathbf{A}$  with respect to  $\mathbf{B}$ , then  $\lambda$  is also an eigenvalue of  $\mathbf{B}^{-1}\mathbf{A}$ , and the converse is also true. Similarly we can show that if  $\lambda$  is an eigenvalue of  $\mathbf{A}$  with respect to  $\mathbf{B}$ , then  $\lambda$  is also an eigenvalue of  $\mathbf{AB}^{-1}$  and of  $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$ , where  $\mathbf{B}^{-1/2}$  is symmetric. The following result summarizes these observations.

**Lemma B.4.1.** *Eigenvalues of  $\mathbf{AB}^{-1}$ ,  $\mathbf{B}^{-1}\mathbf{A}$ , and  $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$  are the same.*

Let  $\lambda_1 \geq \dots \geq \lambda_n$  be the generalized eigenvalues of  $\mathbf{A}$  with respect to  $\mathbf{B}$ . By Lemma B.4.1,  $\lambda_1, \dots, \lambda_p$  are also the eigenvalues of  $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$ . By the Spectral Decomposition Theorem, we have

$$\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2} = \sum_{j=1}^n \lambda_j \mathbf{u}_j \mathbf{u}_j^T, \text{ and } \sum_{j=1}^n \mathbf{u}_j \mathbf{u}_j^T = \mathbf{I},$$

where  $\{\mathbf{u}_j\}$  are the orthonormal eigenvectors of the matrix  $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$  corresponding to the eigenvalues  $\{\lambda_j\}$ . We now write down an analog of Theorem B.2.3 for the generalized eigenvalues.

**Theorem B.4.1.** *Let  $\mathbf{A}$  and  $\mathbf{B}$  be symmetric  $n \times n$  matrices, and assume that  $\mathbf{B}$  is positive definite. Let  $\lambda_1 \geq \dots \geq \lambda_n$  be the eigenvalues of  $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$  with the corresponding orthonormal eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_n$ . It is understood that  $\mathbf{x}, \mathbf{x}_1, \dots$  written below are in  $\mathbb{R}^n$ . Then the following hold:*

- (a)  $\sup_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Bx}} = \lambda_1$ , and this supremum is attained at  $\mathbf{x} = \mathbf{B}^{-1/2} \mathbf{u}_1$ .
- (b)  $\inf_{\mathbf{x} \neq 0} \frac{\mathbf{x}^T \mathbf{Ax}}{\mathbf{x}^T \mathbf{Bx}} = \lambda_n$ , and this infimum is attained at  $\mathbf{x} = \mathbf{B}^{-1/2} \mathbf{u}_n$ .
- (c) For  $1 < m \leq n$ , we have

$$\begin{aligned} & \sup \left\{ \frac{\mathbf{x}_1^T \mathbf{Ax}_1}{\mathbf{x}_1^T \mathbf{Bx}_1} + \dots + \frac{\mathbf{x}_m^T \mathbf{Ax}_m}{\mathbf{x}_m^T \mathbf{Bx}_m} : \mathbf{x}_i^T \mathbf{Bx}_j = 0, 1 \leq i \neq j \leq m, \mathbf{x}_i \neq 0, i = 1, \dots, m \right\} \\ &= \lambda_1 + \dots + \lambda_m, \end{aligned}$$

and this supremum is attained at  $\mathbf{x}_1 = \mathbf{B}^{-1/2} \mathbf{u}_1, \dots, \mathbf{x}_m = \mathbf{B}^{-1/2} \mathbf{u}_m$ .

## B.5 Matrix Derivatives

In many cases one needs to differentiate the quadratic form or the trace or the determinant of a matrix. There are a number of useful formulas for such purposes. For any function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the column vector of its first order partial derivatives and the  $n \times n$  matrix of second order partial derivatives (also known as the Hessian) will be denoted by  $\frac{\partial}{\partial \mathbf{x}} f(\mathbf{x})$  and

$\frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}} f(\mathbf{x})$ , respectively. Let  $\mathbf{A}$  be a symmetric matrix of order  $n \times n$ , and let  $l(\mathbf{x}) = \mathbf{A}\mathbf{x}$  and  $q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ ,  $\mathbf{x} \in \mathbb{R}^n$ .

**Lemma B.5.1.** *Let  $l(\mathbf{x})$  and  $q(\mathbf{x})$  be defined as above. Then,*

$$\frac{\partial}{\partial \mathbf{x}} l(\mathbf{x}) = \mathbf{A}, \quad \frac{\partial}{\partial \mathbf{x}} q(\mathbf{x}) = 2\mathbf{A}\mathbf{x}, \quad \frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}} q(\mathbf{x}) = 2\mathbf{A}.$$

For the result given below, it is assumed that the square matrix  $\mathbf{A}(x)$  of order  $n$  is a function of a real variable  $x$  and element  $(i,j)$  of  $\mathbf{A}(x)$  is  $a_{ij}(x)$ . Let  $\dot{\mathbf{A}}(x)$  denote the matrix obtained by differentiating all elements  $\mathbf{A}(x)$  with respect to  $x$  (ie,  $\dot{\mathbf{A}}(x) = ((\dot{a}_{ij}(x)))$ , where  $\dot{a}_{ij}(x) = da_{ij}(x)/dx$ ). The following result presents expressions for the derivatives of  $\mathbf{A}^{-1}(x)$  and  $|\mathbf{A}(x)|$  with respect to  $x$ .

**Lemma B.5.2.** *Let  $\mathbf{A}(x)$  be an  $n \times n$  symmetric nonsingular matrix whose elements depend on a real variable  $x$ , and let  $\dot{\mathbf{A}}(x)$  be the matrix obtained by differentiating the elements of  $\mathbf{A}(x)$  with respect to  $x$ . Then,*

$$\begin{aligned} \frac{d}{dx} \mathbf{A}^{-1}(x) &= -\mathbf{A}^{-1}(x) \dot{\mathbf{A}}(x) \mathbf{A}^{-1}(x), \\ \frac{d}{dx} |\mathbf{A}(x)| &= |\mathbf{A}(x)| \text{trace}(\dot{\mathbf{A}}(x) \mathbf{A}^{-1}(x)). \end{aligned}$$

It is fairly easy to check the first equality. Since  $\mathbf{I} = \mathbf{A}^{-1}(x)\mathbf{A}(x)$ , differentiating both sides with respect to  $x$ , we have

$$\begin{aligned} \mathbf{0} &= \frac{d}{dx} \mathbf{I} = \frac{d}{dx} \mathbf{A}^{-1}(x) \mathbf{A}(x) \\ &= \left[ \frac{d}{dx} \mathbf{A}^{-1}(x) \right] \mathbf{A}(x) + \mathbf{A}^{-1}(x) \left[ \frac{d}{dx} \mathbf{A}(x) \right] \\ &= \left[ \frac{d}{dx} \mathbf{A}^{-1}(x) \right] \mathbf{A}(x) + \mathbf{A}(x)^{-1} \dot{\mathbf{A}}(x). \end{aligned}$$

Postmultiply by  $\mathbf{A}^{-1}(x)$  on both sides to get the result.

The proof of the second equality is a bit more involved than that of the first. We outline a proof via spectral decomposition of  $\mathbf{A}(x) = \sum_{j=1}^n \lambda_j(x) \mathbf{u}_j(x) \mathbf{u}_j^T(x)$ , where for each  $x$ ,  $\{\mathbf{u}_j(x): j = 1, \dots, n\}$  are orthonormal. Denote the derivatives of  $\lambda_j(x)$  and  $\mathbf{u}_j(x)$  by  $\dot{\lambda}_j(x)$  and  $\dot{\mathbf{e}}_j(x)$ , respectively. Since  $|\mathbf{A}(x)| = \lambda_1(x) \cdots \lambda_n(x)$ , we have

$$\begin{aligned} \frac{d}{dx} |\mathbf{A}(x)| &= \frac{d}{dx} [\lambda_1(x) \cdots \lambda_n(x)] = \sum_{j=1}^n [\lambda_1(x) \cdots \lambda_n(x)] \left[ \dot{\lambda}_j(x) / \lambda_j(x) \right] \\ &= |\mathbf{A}(x)| \sum_{j=1}^n \dot{\lambda}_j(x) / \lambda_j(x). \end{aligned}$$

The second equality would hold if we can establish that

$$\text{trace}(\dot{\mathbf{A}}(x) \mathbf{A}^{-1}(x)) = \sum_{j=1}^n \dot{\lambda}_j(x) / \lambda_j(x).$$

Since  $\mathbf{A}(x) = \sum_{j=1}^n \lambda_j(x) \mathbf{u}_j(x) \mathbf{u}_j^T(x)$ , differentiating both sides with respect to  $x$  we have

$$\begin{aligned}\dot{\mathbf{A}}(x) &= \sum_{j=1}^n \dot{\lambda}_j(x) \mathbf{u}_j(x) \mathbf{u}_j^T(x) + \sum_{j=1}^n \lambda_j(x) \dot{\mathbf{u}}_j(x) \mathbf{u}_j^T(x) + \sum_{j=1}^n \lambda_j(x) \mathbf{u}_j(x) \dot{\mathbf{u}}_j^T(x) \\ &:= \mathbf{B}_1(x) + \mathbf{B}_2(x) + \mathbf{B}_3(x).\end{aligned}$$

The result now follows nothing that

$$\begin{aligned}\mathbf{A}^{-1}(x) &= \sum_{j=1}^n \lambda_j(x)^{-1} \mathbf{u}_j(x) \mathbf{u}_j^T(x), \\ \text{trace}(\mathbf{B}_1(x)\mathbf{A}(x)^{-1}) &= \sum_{j=1}^n \dot{\lambda}_j(x)/\lambda_j(x), \text{ and} \\ \text{trace}(\mathbf{B}_2(x)\mathbf{A}(x)^{-1}) &= \text{trace}(\mathbf{B}_3(x)\mathbf{A}(x)^{-1}) = \sum_{j=1}^n \mathbf{u}_j^T(x) \dot{\mathbf{u}}_j(x) = 0,\end{aligned}$$

where the last step is justified as  $\mathbf{u}_j^T(x) \dot{\mathbf{u}}_j(x) = 0$  for all  $j$ , which can be verified by differentiating both sides of the identity  $\mathbf{u}_j^T(x) \mathbf{u}_j(x) = 1$  with respect to  $x$ .

## B.6 Orthogonal Projection

For a matrix  $\mathbf{A}$  of order  $n \times m$ , we denote its column space  $\{\mathbf{Ax}: \mathbf{x} \in \mathbb{R}^m\}$  by  $\mathcal{M}(\mathbf{A})$ . The orthogonal complement of  $\mathcal{M}(\mathbf{A})$ , denoted by  $\mathcal{M}(\mathbf{A})^\perp$ , is the set  $\{\mathbf{y} \in \mathbb{R}^n: \mathbf{y}^T \mathbf{u} = 0 \text{ for any } \mathbf{u} \in \mathcal{M}(\mathbf{A})\}$ . A square matrix  $\mathbf{A}$  of order  $n$  is called idempotent if  $\mathbf{A}^2 = \mathbf{A}$ . A symmetric matrix  $\mathbf{A}$  is called a (orthogonal) projection matrix if it is symmetric and idempotent. It is fairly easy to see that if  $\mathbf{A}$  is a projection matrix, then so is  $\mathbf{I} - \mathbf{A}$ . Since  $(\mathbf{I} - \mathbf{A})\mathbf{A} = \mathbf{A} - \mathbf{A}^2 = \mathbf{A} - \mathbf{A} = \mathbf{0}$ , it follows that  $\mathcal{M}(\mathbf{A})^\perp = \mathcal{M}(\mathbf{I} - \mathbf{A})$ . If  $\lambda$  is an eigenvalue of  $\mathbf{A}$ , then  $\lambda^2$  is an eigenvalue of  $\mathbf{A}^2$ . Since  $\mathbf{A} = \mathbf{A}^2$ , we have  $\lambda = \lambda^2$  and thus  $\lambda = 0$  or 1. The following lists a few important properties of a projection matrix.

**Theorem B.6.1.** *Let  $\mathbf{A}$  be a  $n \times n$  projection matrix. Let  $\mathcal{M}(\mathbf{A}) = \{\mathbf{Ax}: \mathbf{x} \in \mathbb{R}^n\}$  be the column space of  $\mathbf{A}$ . The following hold:*

- (a)  $\mathbf{I} - \mathbf{A}$  is a projection matrix.
- (b) All the eigenvalues of  $\mathbf{A}$  are either 0 or 1.
- (c)  $\text{trace}(\mathbf{A}) = \text{rank}(\mathbf{A})$ .
- (d)  $\mathcal{M}(\mathbf{A})^\perp = \mathcal{M}(\mathbf{I} - \mathbf{A})$ .
- (e) If  $\mathbf{B}$  is an  $n \times n$  projection matrix and  $\mathcal{M}(\mathbf{B}) \subset \mathcal{M}(\mathbf{A})$ , then  $\mathbf{A} - \mathbf{B}$  is a projection on  $\mathcal{M}(\mathbf{A}) \cap \mathcal{M}(\mathbf{B})^\perp$  and  $\mathbf{AB} = \mathbf{B}$  (ie,  $(\mathbf{I} - \mathbf{A})\mathbf{B} = \mathbf{0}$ ).

Suppose that  $\mathbf{A}$  is of order  $n \times m$ , with rank  $m \leq n$ . Given a vector  $\mathbf{y} \in \mathbb{R}^n$ , how do we find a vector in  $\mathcal{M}(\mathbf{A})$  that is closest to  $\mathbf{y}$ ? Clearly, this is equivalent to minimizing  $\|\mathbf{y} - \mathbf{Ax}\|^2$  with respect to  $\mathbf{x} \in \mathbb{R}^m$ , and if a minimum is attained at  $\mathbf{x} = \mathbf{x}^*$ , then  $\hat{\mathbf{y}} = \mathbf{Ax}^*$  is the

element in  $\mathcal{M}(\mathbf{A})$  that is closest to  $\mathbf{y}$ . It is fairly easy to verify that  $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$  and  $\hat{\mathbf{y}} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{Q}_\mathbf{A} \mathbf{y}$ , where  $\mathbf{Q}_\mathbf{A} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ . It is also easy to check that  $\mathbf{Q}_\mathbf{A}$  is a projection matrix and  $\mathcal{M}(\mathbf{A}) = \{\mathbf{Q}_\mathbf{A} \mathbf{u} : \mathbf{u} \in \mathbb{R}^n\}$ .

## B.7 Distribution of Quadratic Forms

In this section, all the matrices  $\mathbf{A}$ ,  $\mathbf{A}_1$ , etc., associated with quadratic forms of  $\mathbf{Y}$ , where  $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \mathbf{I})$ , are assumed to be symmetric of order  $n \times n$ . If  $\mathbf{A}$  is a projection matrix of rank  $p \leq n$ , then it has exactly  $p$  eigenvalues equal to 1 (and the rest are 0), and let  $\mathbf{u}_1, \dots, \mathbf{u}_p$  be the corresponding orthonormal eigenvectors. Then

$$\|\mathbf{AY}\|^2 = \mathbf{Y}^T \mathbf{AY} = \sum_{j=1}^p (\mathbf{u}_j^T \mathbf{Y})^2 := \sum_{j=1}^p W_j^2.$$

Since  $\{\mathbf{u}_j\}$  are orthonormal,  $\{W_j = \mathbf{u}_j^T \mathbf{Y}, j = 1, \dots, p\}$  are independent with  $W_j \sim N(\mathbf{u}_j^T \boldsymbol{\mu}, 1)$ . Results from [Section 2.2.9](#) in [Chapter 2](#) tell us  $\|\mathbf{AY}\|^2 = \sum_{j=1}^p W_j^2 \sim \chi_p^2(\delta^2)$ , where  $\delta^2 = (1/2) \sum_{j=1}^p (\mathbf{u}_j^T \boldsymbol{\mu})^2 = (1/2) \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} = (1/2) \|\mathbf{A} \boldsymbol{\mu}\|^2$ .

It turns out that the converse is also true, that is, if  $\mathbf{Y}^T \mathbf{AY} \sim \chi_p^2(\delta^2)$ , then  $\mathbf{A}$  must be a projection matrix of rank  $p$ . In order to see this, let us assume that  $\mathbf{A}$  has rank  $q$  and its nonzero eigenvalues and the corresponding normalized eigenvectors are  $\lambda_1, \dots, \lambda_q$  and  $\mathbf{u}_1, \dots, \mathbf{u}_q$ , respectively. Then

$$\mathbf{Y}^T \mathbf{AY} = \sum_{j=1}^q \lambda_j W_j^2,$$

where  $\{W_j = \mathbf{u}_j^T \mathbf{Y}, j = 1, \dots, q\}$  are independent with  $W_j \sim N(\mathbf{u}_j^T \boldsymbol{\mu}, 1)$ . Since  $W_j^2 \sim \chi_1^2(\delta_j^2)$  with  $\delta_j^2 = (1/2) (\mathbf{u}_j^T \boldsymbol{\mu})^2$ , and  $\{W_j : j = 1, \dots, q\}$  are independent, the characteristic function (cf) of  $\mathbf{Y}^T \mathbf{AY}$  is the product of the cf's of  $\lambda_j W_j^2, j = 1, \dots, q$ . And this product of cf's must be equal the cf of  $\chi_p^2(\delta^2)$ , since  $\mathbf{Y}^T \mathbf{AY} \sim \chi_p^2(\delta^2)$  by assumption. An examination of this equality of the characteristic functions shows (details not given here) that  $p$  must be equal to  $q$  and  $\lambda_1 = \dots = \lambda_q = 1$ . This proves that  $\mathbf{A}$  is a projection matrix of rank  $p$  and we have the following result.

**Lemma B.7.1.** *If  $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \mathbf{I})$ , then  $\mathbf{Y}^T \mathbf{AY} \sim \chi_p^2(\delta^2)$  if and only if  $\mathbf{A}$  is a projection matrix of rank  $p$ .*

We now discuss a more general result. Let  $\mathbf{A}_1$  and  $\mathbf{A}_2$  be two projection matrices of ranks  $p_1$  and  $p_2$ , respectively, and  $\mathbf{A}_1 \mathbf{A}_2 = 0$ . We can therefore find orthonormal vectors  $\{\mathbf{u}_j : j = 1, \dots, p_1 + p_2\}$  such that

$$\mathbf{A}_1 = \sum_{j=1}^{p_1} \mathbf{u}_j \mathbf{u}_j^T \text{ and } \mathbf{A}_2 = \sum_{j=p_1+1}^{p_1+p_2} \mathbf{u}_j \mathbf{u}_j^T.$$

Since  $\{W_j = \mathbf{u}_j^T \mathbf{Y}, j = 1, \dots, p_1 + p_2\}$  are independent with  $W_j \sim N(\mathbf{u}_j^T \boldsymbol{\mu}, 1)$ , we can conclude that

$$\begin{aligned}\|\mathbf{A}_1 \mathbf{Y}\|^2 &= \mathbf{Y}^T \mathbf{A}_1 \mathbf{Y} = \sum_{j=1}^{p_1} W_j^2 \sim \chi_{p_1}^2(\delta_1^2), \\ \|\mathbf{A}_2 \mathbf{Y}\|^2 &= \mathbf{Y}^T \mathbf{A}_2 \mathbf{Y} = \sum_{j=p_1+1}^{p_1+p_2} W_j^2 \sim \chi_{p_2}^2(\delta_2^2),\end{aligned}$$

and that  $\mathbf{Y}^T \mathbf{A}_1 \mathbf{Y}$  and  $\mathbf{Y}^T \mathbf{A}_2 \mathbf{Y}$  are independent, where  $\delta_1^2 = (1/2)\boldsymbol{\mu}^T \mathbf{A}_1 \boldsymbol{\mu}$  and  $\delta_2^2 = (1/2)\boldsymbol{\mu}^T \mathbf{A}_2 \boldsymbol{\mu}$ . Moreover,  $\mathbf{Y}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{Y} \sim \chi_{p_1+p_2}^2(\delta_1^2 + \delta_2^2)$ .

It turns out that a converse of this is also true as given in the following result. The proofs of the next two results use ideas similar to the ones given above and details can be found in Rao [66].

**Lemma B.7.2.** *Let  $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$  and assume that*

- (i)  $\mathbf{Y}^T \mathbf{A} \mathbf{Y} \sim \chi_p^2(\delta^2)$ ,
- (ii)  $\mathbf{Y}^T \mathbf{A}_1 \mathbf{Y} \sim \chi_{p_1}^2(\delta_1^2)$ , and
- (iii)  $P[\mathbf{Y}^T \mathbf{A}_2 \mathbf{Y} \geq 0] = 1$ .

*Then  $\mathbf{A}_2$  is a projection matrix of rank  $p - p_1$  and  $\mathbf{A}_1 \mathbf{A}_2 = \mathbf{0}$ . Consequently,  $\mathbf{Y}^T \mathbf{A}_2 \mathbf{Y} \sim \chi_{p-p_1}^2(\delta^2 - \delta_1^2)$  and  $\mathbf{Y}^T \mathbf{A}_1 \mathbf{Y}$  is independent of  $\mathbf{Y}^T \mathbf{A}_2 \mathbf{Y}$ .*

**Lemma B.7.3.** *Assume that  $\mathbf{Y}^T \mathbf{A}_1 \mathbf{Y} \sim \chi_{p_1}^2(\delta_1^2), \dots, \mathbf{Y}^T \mathbf{A}_r \mathbf{Y} \sim \chi_{p_r}^2(\delta_r^2)$ . Then a necessary and sufficient condition that  $\mathbf{Y}^T \mathbf{A}_1 \mathbf{Y}, \dots, \mathbf{Y}^T \mathbf{A}_r \mathbf{Y}$  are independent is that  $\mathbf{A}_j \mathbf{A}_k = \mathbf{0}$  for all  $j \neq k$ , in which case,  $\mathbf{Y}^T (\mathbf{A}_1 + \dots + \mathbf{A}_r) \mathbf{Y} \sim \chi_{p_1+\dots+p_r}^2(\delta_1^2 + \dots + \delta_r^2)$ .*

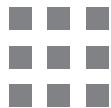


# Bibliography

1. Ferguson T. *Mathematical Statistics*. New York: Academic Press; 1967.
2. Lehmann E. *Testing Statistical Hypotheses*. New York: John Wiley; 1959.
3. Lehmann E. *Theory of Point Estimation*. New York: John Wiley; 1983.
4. Hájek J, Šidák Z. *Theory of Rank Tests*. New York: Academic Press; 1967.
5. Gosset W. On the probable error of a mean. *Biometrika* 1908;6:1–25.
6. Kolmogorov A. *Foundations of the Theory of Probability*. New York: Chelsea Publishing Company; 1933 (English Translation, 1950).
7. Hoeffding W. Probability inequalities for sums of bounded random variables. *J Am Statist Assoc* 1963;58:13–30.
8. Uspensky J. *Introduction to Mathematical Probability*. New York and London: McGraw Hill; 1937.
9. Doob J. *Stochastic Processes*. New York: John Wiley; 1953.
10. Hoeffding W, Robbins H. The central limit theorem for dependent random variables. *Duke Math J* 1948;15:773–80.
11. Stein C. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In: Proc Third Berkeley Symp Math Statist Prob, vol. 1. Univ of Calif Press, CA; 1956. p. 197–206.
12. Hodges J, Lehmann E. Some problems in minimax point estimation. *Ann Math Statist* 1950;21:182–97.
13. Karlin S. Admissibility for estimation with quadratic loss. *Ann Math Statist* 1958;29:406–36.
14. Stein C. The admissibility of Pitman's estimator for a single location parameter. *Ann Math Statist* 1959;30:970–9.
15. Rao C. Information and accuracy attainable in the estimation of statistical parameters. *Bull Calc Math Soc* 1945;37:81–91.
16. Blackwell D. Conditional expectation and unbiased sequential estimation. *Ann Math Statist* 1947;18:105–10.
17. Basu D. On statistics independent of a complete sufficient statistic. *Sankhya* 1955;15:377–80.
18. Cramér H. *Mathematical Methods of Statistics*. Princeton, NJ: Princeton University Press; 1946.
19. Pitman E. The estimation of location and scale parameters of a continuous population of any given form. *Biometrika* 1939;30:391–421.
20. Fisher R. On the mathematical foundations of theoretical statistics. *Phil Trans R Soc Lond Ser A* 1921;222:309–68.
21. Rao C. Criteria of estimation in large samples. *Sankhya* 1963;25:189–206.
22. Neyman J, Pearson E. On the problem of the most efficient tests of statistical hypotheses. *Phil Trans R Soc Lond Ser A* 1933;231:289–337.
23. Neyman J, Pearson E. Contributions to the theory of testing statistical hypotheses. I. Unbiased critical regions of type A and type A<sub>1</sub>. *Stat Res Mem* 1936;1:1–37.
24. Wald A. *Sequential Analysis*. New York: John Wiley; 1947.

25. Bahadur R. A note on the fundamental identity of sequential analysis. *Ann Math Statist* 1958;29:534–43.
26. Neyman J. Outline of a theory of statistical estimation based on the classical theory of probability. *Phil Trans R Soc* 1937;235:333–80.
27. Wald A. Note on the consistency of the maximum likelihood estimate. *Ann Math Statist* 1949;20:595–601.
28. Neyman J, Scott E. Consistent estimates based on partially consistent observations. *Econometrika* 1948;16:1–32.
29. Basu D. An inconsistency of the method of maximum likelihood. *Ann Math Statist* 1955;26:144–5.
30. Ferguson T. An inconsistent maximum likelihood estimate. *J Am Statist Assoc* 1982;77:831–4.
31. LeCam L. On some asymptotic properties of maximum likelihood estimates and related Bayes estimates. *Univ Calif Publ Statist* 1953;1:277–330.
32. Chernoff H, Lehmann E. The use of maximum likelihood estimates in  $\chi^2$  tests for goodness of fit. *Ann Math Statist* 1954;25:579–86.
33. Hoeffding W. A class of statistics with asymptotic normal distribution. *Ann Math Statist* 1948;19:293–325.
34. Roussas G. *Contiguity of Probability Measures*. Cambridge: Cambridge University Press; 1972.
35. Fisher R, Yates F. *Statistical Tables for Biological, Agricultural and Medical Research*. Edinburgh-London: Oliver and Boyd; 1938.
36. Hoeffding W. “Optimum” nonparametric tests. In: Proc Second Berkeley Symp Math Statist Prob, vol. 1. Univ Calif Press, CA; 1950. p. 83–92.
37. Terry M. Some rank order tests which are most powerful against specific parametric alternatives. *Ann Math Statist* 1952;23:346–66.
38. Chenoff H, Savage I. Asymptotic normality and efficiency of certain nonparametric test statistics. *Ann Math Statist* 1958;29:972–94.
39. Dvoretzky A, Kiefer J, Wolfowitz J. Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *Ann Math Statist* 1956;27:642–69.
40. Doob J. Heuristic approach to the Kolmogorov-Smirnov theorems. *Ann Math Statist* 1949;20:393–403.
41. Freedman D. *Brownian Motion and Diffusion*. New York: Springer; 1983.
42. Breiman L. *Probability*. Reading, MA: Addison-Wesley; 1968.
43. Billingsley P. *Convergence of Probability Measures*. New York: John Wiley; 1968.
44. Wiener N. Un probleme de probabilités énonmbrables. *Bull Soc Math de France* 1924;52:569–78.
45. Rosenblatt M. Remarks on some nonparametric estimates of a density function. *Ann Math Statist* 1956;27:832–7.
46. Parzen E. On estimation of a probability density and mode. *Ann Math Statist* 1962;33:1065–76.
47. Nadaraya E. On nonparametric estimates of density functions and regression curves. *Theor Prob Appl* 1965;10:186–90.
48. Bhattacharya P. Estimation of a probability density function and its derivatives. *Sankhya Ser A* 1967;29:373–82.
49. Silverman B. *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall; 1968.
50. Fan J, Gijbels I. *Local Polynomial Modelling and its Applications*. London: Chapman-Hall; 1996.
51. Stone C. On asymptotically optimal window selection rule for kernel density estimates. *Ann Statist* 1984;12:1285–97.

52. Kaplan E, Meier P. Nonparametric estimation from incomplete observations. *J Am Statist Assoc* 1958;53:457–81.
53. Efron B. The two-sample problem with censored data. In: Proc Fifth Berkeley Symp Math Statist Prob, vol. 4. University of California Press, CA; 1967. p. 831–53.
54. von Mises R. On the asymptotic distribution of differentiable statistical functionals. *Ann Math Statist* 1947;18:309–48.
55. Serfling R. *Approximation Theorems in Mathematical Statistics*. New York: John Wiley; 1980.
56. Huber P. *Robust Statistics*. New York: John Wiley; 1981.
57. Fernholz L. *von Mises Calculus for Statistical Functionals*. New York: Springer-Verlag; 1983. Lecture Notes in Statistics.
58. Bahadur R. A note on quantiles in large samples. *Ann Math Statist* 1966;37:577–80.
59. Efron B, Tibshirani R.J. *An Introduction to the Bootstrap*. New York: Chapman and Hall; 1993.
60. Searle SR. *Linear Models*. New York: John Wiley; 1971.
61. Mardia K, Kent J, Bibby J. *Multivariate Analysis*. New York: Academic Press; 1979.
62. Box G, Jenkins G, Reinsel G. *Time Series Analysis: Forecasting and Control*. New Jersey: Prentice-Hall; 1994.
63. Brockwell P, Davis R. *Time Series: Theory and Methods*. New York: Springer; 1991.
64. Gikhman I, Skorokhod A. *Introduction to the Theory of Random Processes*. New York: Dover; 1996.
65. Royden H. *Real Analysis*. New York: McMillan; 1968.
66. Rao C. *Linear Statistical Inference and its Applications*. New Delhi: Wiley Eastern; 1973.
67. Hájek J. A characterization of limiting distributions of regular estimates. *Z Wahr Verw Geb* 1970;14:323–30.
68. Inagaki N. On the limiting distribution of a sequence of estimators with uniformity property. *Ann Inst Stat Math* 1970;22:113–97.
69. Pitman E. *Notes on Nonparametric Statistical Inference*. New York (Mimeographed): Columbia Univ; 1948.



# Index

## A

Acceptance region, 126  
*ACF*. *See Autocorrelation function (ACF)*  
Additive model, 314  
Akaike information criterion (AIC)  
heuristic derivation, 349–352  
predictive likelihood, 349  
uniform integrability, 350–351  
Almost sure convergence, 189–191  
Alternative hypothesis  
hypothesis testing, 71, 125  
null hypothesis, 77–78  
simple hypothesis, 126  
two-decision problem, 77  
Analysis of covariance  
(ANCOVA)  
application of, 344–345  
 $\beta$  and  $\gamma$  estimation of, 340–341  
 $\gamma$  and  $\beta$  inference for, 342–344  
residual sum of squares, 342  
Analysis of variance (ANOVA), 309, 311,  
318–319  
Ancillary statistic, 96, 110, 153  
ANCOVA. *See Analysis of covariance (ANCOVA)*  
ANOVA. *See Analysis of variance (ANOVA)*  
ARIMA. *See Integrated autoregressive-moving average (ARIMA)*  
ARMA. *See Autoregressive-moving average (ARMA)*  
ARMA( $p, q$ ) models. *See*  
Autoregressive-moving average  
(ARMA( $p, q$ )) models  
 $AR(p)$  models. *See Autoregressive process*  
( $AR(p)$ ) models  
Asymptotic bias (ABias), 265–266, 299–301  
Asymptotic distributions, 191–192, 408–410  
 $\alpha$ -trimmed mean, 288–289  
delta method, 278, 280–281

*L*-estimator, 286–287  
*M*-estimator, 291–292  
Asymptotic efficacy, 116, 225  
Asymptotic normality, 181, 258, 261  
Asymptotic properties, MLEs  
consistency and asymptotic normality,  
179–194  
frequency  $\chi^2$ , 202–208  
independent nonidentically distributed  
data, 201–202  
likelihood ratio test, 194–201  
Asymptotic relative efficiency (ARE), 224  
Asymptotic variance (AVar)  
asymptotic bias, 265–266, 299–300  
consistent estimators, 182  
estimation of, 300–301  
Augmented matrix, 340, 369  
Autocorrelation function (ACF), 434,  
442–450  
autocovariance functions, 445–448  
backshift operator, notation, 450  
infinite sequence, 444  
linear stationary series, 443, 447  
mean, 442–445  
PACF, 449–450  
Autocovariance function, 434, 437–440  
Autoregressive-moving average (ARMA)  
ACF, 449–450  
appropriate selection, 472–473  
 $AR(p)$  models, 458–459  
diagnostics, 448–449  
forecasting, 463–464  
 $MA(q)$  models, 459  
PACF, 449–450  
Autoregressive-moving average  
(ARMA( $p, q$ )) models  
invertibility, 459–460  
nonredundancy, 459–460

Autoregressive-moving average  
 (ARMA( $p, q$ ) models (*Continued*)  
 parameter estimation, 471–472  
 redundancy issue, 458  
 stationarity, 459–460  
 technical issues, 457  
 Autoregressive process (AR( $p$ ) models  
 condition for, 458–459  
 forecasting, 461  
 parameter estimation, 467–469  
 stationarity, 436  
 Average sampling number (ASN) function, 162, 166–168

**B**

Bahadur sample quantile, 283–284  
 Balanced two-factor ANOVA, 323–324, 331  
 Bandwidth, 256–257  
 Bayes formula, 6, 80–81  
 Bayesian estimation, 112–114  
 Bayesian information criterion (BIC)  
 discrete variable, 352  
 Laplace approximation, 353  
 positive definite matrix, 352–353  
 Bayes principle, 79  
 Bayes' rule, 80–81, 415  
 Behavioral decision rule, 72, 126  
 Bernoulli parameter, 72  
 Bernoulli ( $p$ ) rv, 25–26  
 Bernstein's inequality, 63  
 Best asymptotically normal (BAN), 182, 193–194  
 Best linear unbiased estimator (BLUE), 318, 364  
 Best linear unbiased prediction (BLUP), 364–365  
 Beta distributions, 34–36  
 Beta functions, 33–34  
 Bias-squared terms, 358  
 Bivariate distribution, 218, 255  
 Bivariate normal, 170  
 Block diagonal matrix, 340–341  
 Bochner's theorem, 440  
 Bonferroni method, 337

confidence intervals, 397  
 one-sample inference, 390  
 two-sample problem, 395  
 Bootstrap resampling method, 298–303  
 asymptotic bias, 300–301  
 asymptotic variance, 300–301  
 heuristic justification, 301–303  
 Borel-Cantelli Lemma, 57, 241, 259  
 Bounded variation, 259–260, 494  
 Box-Cox transformation, 433

**C**

Calculus, fundamental theorem of, 95  
 Canonical correlation analysis, 420–425  
 cross-classified data, 422–423  
 sample estimates, 421  
 social mobility data, 423  
 technical notes, 424–425  
 test for, 421–422  
 Cauchy distribution, 45–46, 131, 285–286  
 Cauchy-Schwarz inequality, 12, 493, 504  
 Causal time series, 454–460  
 invertible representation, 455–456  
 moving average models, 457–458  
 technical issues, 456–457  
 usefulness of, 456  
 Censored data, 255, 273, 275  
 Censoring time, 273  
 Central Limit Theorem (CLT)  
 asymptotic normality, 66–67  
 characteristic function, 496–500  
 Lindeberg-Lévy, 56  
 multivariate, 385–386, 500  
 Change of variable, 495  
 Characteristic function, 496–500  
 continuity theorem, 498–500  
 Helly-Bray theorem, 496, 498  
 properties of, 497–498  
 special cases, 497  
 uniqueness of, 498  
 Chernoff-Savage approach, 238–239  
 Chi-squared distribution, 40–41, 390–391  
 Classification function  
 Bayes' rule, 417  
 Fisher's method, 418–420

- method of, 419
- probability of, 417–418
- CLT. *See* Central Limit Theorem (CLT)
- Coefficient of determination, 318–319
- Communality, 411
- Composite hypothesis, 131–132
- Conditional distribution, 147, 154
- Conditional expectation, 300–301
- Conditional probability, 347
- Conditional tests, 149–152
- Confidence coefficient, 169–171
- Confidence ellipsoid
  - one-sample inference, 389
  - two-sample problem, 394
- Confidence intervals, 170, 328
  - Bonferroni method, 390, 395
  - confidence coefficient, 410
  - Scheffé method, 389, 395
  - simultaneous, 397, 401–402
- Confidence sets
  - inverting acceptance regions, 171–173
  - pivotal functions, 169–171
- Conjugate gradient method, 360
- Consistency and asymptotic normality
  - almost sure convergence, 189–191
  - efficiency, 191–193
  - multinomial distribution, parameters in, 193–194
- Consistent estimators, 182
- Continuation region, 164
- Continuity theorem, 498–500
  - Cramér-Wold device, 499
  - Helly-Bray theorem, 498
  - Khinchine weak law of large numbers, 498
  - Lindeberg-Lévy theorem, 499
  - multivariate central limit theorem, 500
- Continuous distributions
  - beta distributions, 34–36
  - beta functions, 33–34
  - Cauchy distribution, 45–46
  - Chi-square distribution, 40–41
  - exponential distribution, 36
  - exponential family, 49–51
  - F distributions, 42–45
  - gamma distributions, 34–36
  - gamma functions, 33–34
  - multivariate normal distribution, 46–49
  - noncentral  $\chi^2$ , 43–45
  - normal distribution, 36–42
  - sample mean, 41–42
  - sample variance, 41–42
  - $t$  distributions, 42–43
  - uniform distribution, 34
- Continuous mapping theorem, 501
- Convergence modes, 55–60
  - in law, 56
  - in probability, 56
  - in quadratic mean, 56
- Convex function, 492–493
- Correlation coefficient, 12, 16
- Correlation matrix, 412
- Covariance
  - ANCOVA, application of, 344–345
  - $\beta$  and  $\gamma$  estimation of, 340–345
  - $\gamma$  and  $\beta$  inference for, 342–344
  - residual sum of squares, 342
- Covariance matrix, 11
- Covariance stationarity. *See* Strictly stationary
- Cramér-Rao inequality, 97
- Cramér-Rao information inequality, 99
- Cramér-Rao lower bound, 96–98
- Cramér-Wold device, 468–469
- Critical function, 126, 135
- Critical region/rejection region, 126
- Cross-product term, 323, 325
- Cross-validation method, 269, 353–355
- Cumulative distribution function (cdf), 7
- Curve estimation
  - density, 255–260
  - higher dimension, 264–265
  - local polynomials, 265–272
  - nearest neighbor approach, 263–264
  - regression, 260–263
  - survival function and hazard rates, 273–275
- D**
- Decision functions, 71–74. *See also* Optimal decision rules
- Degrees of freedom (df), 317, 319–320

Delta method  
 Bahadur sample quantile, 283–284  
 differentiability condition, 279  
 gross-error  $\delta_\chi$ , 280  
 influence function, 280, 282–285  
 $k$ th central moment, 282  
 mean, 282  
 partial derivatives, 60  
 $p$ -quantile, 283  
 Taylor expansion, 278  
 variance of, 274  
 $V$ -statistic, 277, 284  
 De Moivre-Laplace theorem, 37–40  
 DeMorgan's rules, 2, 58  
 Density estimation  
   bandwidth choice, 269–270  
   cross-validation procedure, 270–271  
   integrated square-error, 269–270  
   leave-one-out, 269–270  
   optimality property, 271  
 Design matrix, 309, 318–319  
 Diagnostics, 448–450  
 Differential conditions  
    $\rho$ -Fréchet differentiability, 279  
   Gâteaux differentiability, 279  
    $\rho$ -Hadamard differentiability, 279  
 Differentiation under integral, 491  
 Discrete distribution, 255–256  
   binomial distribution, 25–26  
   geometric distribution, 27–28  
   hypergeometric distribution, 29–30  
   multinomial distribution, 26–27  
   negative binomial distribution, 28–29  
   poisson distribution, 30–33  
 Discrete Fourier transform, 476–477  
 Disjoint subsets, 125–126  
 Distribution-free property, 242  
 Dominated convergence theorem, 491  
 Double exponential, 358  
 Durbin-Levinson iterative, 451–453  
 Durbin-Levinson recursions, 450–451

**E**  
 Eigenvalues, 205–206, 504–506  
   singular value decomposition, 505  
   spectral decomposition theorem, 505  
 Eigenvectors, 505–506  
 Elementary facts  
   Cauchy-Schwarz inequality, 504  
   Sherman-Morrison formula, 504  
   symmetric matrix  $A$ , 503  
 Empirical distribution function  
   asymptotic distributions, 246–249  
   Brownian motion, 243–245  
   test statistics, 241–243  
   weak convergence, 245  
 Empirical distribution function (edf), 55  
 Equivariance, 106–112  
 Equivariant under location, 77  
 Error probabilities, 163–164  
 Error vector, 315–316  
 Estimation  
    $\beta$  and  $\Sigma$ , 403–404  
   canonical correlation analysis, 421  
   factor analysis, 411  
   MANOVA model, 395–396  
   principal components, 407–408  
   properties of, 404  
   two-factor MANOVA, 400  
   two-sample problem, 393–394  
 Euclidean space, 135  
 Exact distributions, 397–398  
 Expected value, random variable, 495  
 Exponential distributions, 36, 199  
 Exponential family, distribution, 49–51, 98  
 Extracting stationary part, 433–434  

**F**

 Factor analysis, 411–414  
   estimation of, 411  
   maximum likelihood, 412  
   prediction of, 413–414  
   principal, 411–412  
 Factor-effect smodel, 311–313  
 Factorization theorem, 76, 388  
 $F$ -distribution, 42–45, 329  
 Feller's theorem, 170

- Finite expectation, 214  
 Finite sample space, 4–5  
 Fisher-information  
     Cramér-Rao bound, 118  
     exponential families, 98  
     information inequality, 97  
      $M$ -estimator, 290–291  
     Wald's statistics, 211  
 Fisher-Irwin test, 152  
 Fisher's method, 418–420  
 Fitted mean vector, 317, 335  
 $\rho$ -Fréchet differentiability, 279, 294–295  
 Frequency distribution, 207  
 Frequency response function, 481–483  
 Frequency  $x^2$ , 202–208  
 $F$ -statistic, 328–329  
 Fubini theorem, 492  
 Fundamental identity, 167
- G**  
 Gamma distributions, 34–36  
 Gamma functions, 33–34  
 Gâteaux differentiability, 279–280, 297  
 Gaussian stationary, 434–435  
 Gauss-Markov models, 402, 468  
      $\beta$  and  $\sigma^2$  estimation of, 316–317  
     Bonferroni method, 337  
     inference, 328–340  
     linear functions, estimation of, 317–318  
     linear unbiased estimation, 318  
     one-factor balanced ANOVA model, 339  
     prediction intervals, 339–340  
     Scheffé method, 337–338  
     simultaneous, 336–337  
     Tukey method, 338–339  
 Generalized cross-validation (GCV), 346  
 Generalized eigenvalues, 507  
 Gross-error  $\delta_\chi$ , 280
- H**  
 $\rho$ -Hadamard differentiability, 279  
 Hájek projection method, 220–221  
 Hájek-Rényi inequality, 64  
 Hardy-Weinberg formula, 212  
 Hazard function, 275
- Heine-Borel property, 191  
 Helly-Bray theorem, 496, 498  
 Higher dimension  
     bias, 265  
     curse of dimensionality, 265  
     iid observation, 264  
     kernel method, 265  
     optimal bandwidth, 265  
     regression function, 264  
     variance, 265  
 Histogram, 255–256  
 Hoeffding's inequality, 62  
 Hölder's inequality, 165, 493  
 Homogeneity distributions, 199  
 Homogeneity probabilities, 200  
 Hotelling's  $T^2$ -distribution, 384  
 Hypergeometric distribution, 131, 152  
 Hypothesis testing  
     conditional tests, 149–152  
     confidence sets, 169–173  
     empirical distribution function, 241–249  
     generalized Neyman-Pearson lemma, 135–136  
     locally best tests, 140–144  
     MANOVAmode, 396–397  
     one-sample inference, 390  
     one-sided hypotheses, UMP tests for, 131–132  
      $p$ -value, 159–160  
     ranks and order statistics, 213–227  
     rank tests, 227–241  
     sequential probability ratio test, 160–168  
     simple null hypothesis *vs.* simple alternative, 127–131  
     two-factor MANOVA, 401  
     two-sample problem, 394  
     two-sided problems, UMP tests for, 135–136  
     unbiased tests, 133–135
- I**  
 Identifiability condition, 117, 189  
 Independent nonidentically distributed data, 201–202  
 Independent variable, 309–310  
 Influence function, 280, 282–285

- Information inequality  
 Cramér-Rao lower bound, 96–98  
 in multiparameter families, 99–106
- Information lower bound, 102
- Information matrix, 102
- Initial estimator, 185–189
- Instantaneous failure rate, 273
- Integrated autoregressive-moving average (ARIMA)  
 forecasting, 466–467  
 stationary part, 436
- Integrated square error, 258
- Integration theory, 491–492
- Interaction effects, 313–314
- Intraclass correlation, 362–363
- Invariant under location, 77
- Inverse of matrix, 506
- Invertible time series. *See* Causal time series
- J**
- Jackknife, 298–303  
 asymptotic bias, 300–301  
 asymptotic variance, 300–301  
 heuristic justification, 301–303
- Jensen's inequality, 492–493
- Joint distribution, 144, 213–214
- K**
- Kendall's tau statistic, 218
- Kernel estimator, properties of, 256–260
- Kolmogorov's inequality, 63
- $\kappa$ th central moment, 282
- Kullback-Leibler divergence, 346–347
- L**
- Lagrangian multiplier, 326–327
- Lasso method, 357–358
- Law of large number, 496–500
- Least square estimate, 328, 358–359
- Leave-one-out, 269–270
- Lebesgue dominated convergence, 165
- Left-sided null hypothesis *vs.* right-sided alternative, 171
- Lehmann-Scheffé theorem, 103
- $L$ -estimators, 285–289  
 $\alpha$ -trimmed mean, 288–289  
 asymptotic distribution, 286–287  
 Cauchy distribution, 285–286  
 $M$ -estimators, 293  
 score function, 293
- Level of significance, 126
- Likelihood equation, 188
- Likelihood function ( $L$ ), 115, 387, 390–391
- Likelihood ratio statistic, 211
- Likelihood ratio test (LRT), 194–201, 390–392  
 factor A main effects, 401  
 factor B main effects, 401  
 for interaction test, 401
- Lindeberg condition, 262
- Lindeberg-Feller theorem, 67
- Lindeberg-Liapounov theorem, 67
- Linear discriminant rule  
 Fisher's method, 419  
 sample estimates, 416–417
- Linear filtering, 480–483
- Linear functions, 317–318
- Linear models, 402–404, 414  
 $\beta$ , linear restrictions, 325–328  
 covariance analysis, 340–345  
 Gauss-Markov models, 310–315  
 inference, 374–378  
 model selection, 345–355  
 random and mixed-effects, 361–373  
 regression, methods for, 356–361  
 total sum of squares, decomposition of, 318–325
- Linear prediction, 441–442
- Linear regression model, 309, 311
- Linear time series, 440–441
- Ljung-Box test, 448
- Loading vectors, 405
- Local linear estimate, 266–268
- Locally best tests  
 Fisher-information, 141–142  
 locally most powerful, 140  
 logistic distribution, 143  
 multiparameter exponential families, 140  
 random sample, 140–142  
 regularity conditions, 140  
 UMP unbiased tests, 140

- Locally most powerful (LMP), 140, 228
- Local polynomials
- asymptotic bias, 265–267
  - asymptotic variance, 265–267
  - cross-validation method, 269
  - density estimation, 269–271
  - local linear estimate, 266–268
  - regression estimation, 271–272
  - regression function, 265–266
  - regression model, 266–267
- Location-scale family, 170
- Logistic distribution, 143
- Log likelihood ratio, 161–162
- Loss function, 72
- Lower bound
- Cramér-Rao lower bound, 96–98
  - Fisher-information, 98
  - multiparameter families, 99–106
- Lower confidence bound, 169
- M**
- Mahalanobis distance, 385, 414
- $MA(q)$  models. *See* Moving average ( $MA(q)$ ) models
- Mann-Whitney statistic, 217–218
- MANOVA model
- confidence intervals, 397
  - estimation, 395–396
  - hypothesis testing, 396–397
  - one-factor tests, 398–399
  - of test interpretation, 399
  - Wilks' lambda, 397–398
- Marginal distribution, 147
- Markov inequality, 60
- Martingale property, 64
- Matrix algebra
- distribution of quadratic forms, 510–511
  - eigenvalues and eigenvectors, 504–506
  - elementary facts, 503–504
  - generalized eigenvalues, 507
  - matrix derivatives, 507–509
  - orthogonal projection, 509–510
  - symmetric matrix function, 506
- Matrix derivatives, 507–509
- Maximum likelihood estimators (MLEs)
- $M$ -estimator, 289
  - method of, 115
  - normal population sampling, 387
  - variance components, 371–373
  - Wilks' lambda, 398
- Mean, 442–445. *See also* Autocorrelation function
- forecast error, 456
  - influence function, 282
- Mean-square error (MSE)
- bias and variance, 256–257
  - Mallows' criterion, 347–348
  - unbiased estimator, 317
- $M$ -estimator, 289–292. *See also* Maximum likelihood estimators (MLEs)
- asymptotic distribution, 291–292
  - Huber functions, 291
  - $L$ -estimators, 293
  - $L$ -functional, 295
  - minimax problem, 291
  - monotone score function, 291–292
  - score function, 293
- Method of maximum likelihood, 115–118
- Method of minimum  $\chi$ , 119–121
- Method of moments estimators (MOME), 119
- Minimax principle, 79
- Minimax rules, 82–83
- Minimum norm quadratic unbiased estimation (MINQUE), 369
- Minimum risk equivariant (MRE) estimator, 107–108
- Mixed-effects models, 413
- equations, 366–367
  - inference, 374–378
  - variance components, estimation of, 369–371
- Mixed model equations
- assumption of normality, 367–368
  - likelihood function, 367–368
  - motivation for, 367–369
  - one-factor random effects model, 366–367
  - Sherman-Morrison formula, 368–369
- MLEs. *See* Maximum likelihood estimators (MLEs)

MLR. *See* Monotone likelihood ratio (MLR)

Model selection

- AIC and BIC criteria, 348–353
- Akaike's FPE, 347–348
- cross-validation, 353–355
- Mallows' criterion, 347–348

Monotone convergence theorem, 191, 491

Monotone likelihood ratio (MLR), 131–132

Monotone power, 137

Most powerful (MP), 127–128

Moving average ( $MA(q)$ ) models

- forecasting, 461–463
- identifiability of, 459
- nonuniqueness of, 457–458
- parameter estimation, 469–471

MSE. *See* Mean-square error (MSE)

Multinomial coefficient, 26–27, 152

Multinomial distribution, 193–194, 202–203, 212

Multinomial probabilities, 200

Multiple linear regression model, 15

Multivariate analysis

- bonferroni method, 390
- canonical correlation analysis, 420–425
- central limit theorem, 385–386
- classification and discrimination, 414–420
- confidence ellipsoid, 389
- confidence intervals, 389–390
- factor analysis, 411–414
- hypothesis testing, 390
- likelihood ratio test, 390–392
- linear model, 402–404
- mahalanobis distance, 385
- MANOVA model, 395–399
- normality, 386–387
- normal population, sampling, 387
- one-sample inference, 388–392
- principal components analysis, 404–410
- sampling distributions, 387–388
- two-factor MANOVA, 400–402
- two-sample problem, 393–395
- wishart distribution, 383–385

Multivariate normal distribution, 46–49

Multivariate normality, 386–387

**N**

Natural parameter space, 133–135

Nearest neighbor approach

- density estimation, 263
- kernel estimation procedure, 264
- MSE, 263
- regression estimation, 263–264
- second-order smoothness condition, 264

Nested ANOVA model, 334

Newton-Raphson method, 183–189, 193–194

Neyman-Pearson lemma

- corollary, 128–131
- Euclidean space, 135
- likelihood ratio, 129
- MP level, 129–131
- parametric family, 129–131
- uniformly most powerful, 130

Noncentral  $\chi^2$ , 43–45

Noncentrality parameter, 320–322

Nonlinear regression, 311

Nonparametric estimate, 272

Nonparametric models, 70

Nonrandomized decision rule, 126

Nonstationary series, 433

Normal distribution, 36–42

Nuisance parameters

- alternative hypotheses, 147
- exponential family, 146
- joint distribution, 144
- Neyman-structure, 145
- normal distribution, context of, 154–158
- null hypothesis, 144–145
- similarity and completeness, 144–158
- sufficient statistic, 145–146
- three problems, 145

Null hypothesis

- alternative hypothesis, 77
- consistency, 197
- hypothesis testing, 71

**O**

One-factor ANOVA model, 311–312

One random factor, 361

One-to-one transforms, 156, 158

- Operating characteristic (OC) function, 162, 164–166
- Optimal bandwidth, 265
- Optimal decision rules
- Bayes rules, 80–81
  - conditions for admissibility, 83–86
  - estimation problem, 77
  - minimax rules, 82–83
  - suitable ordering of, 78–80
  - two-decision problem, 77–78
- Optimality under unbiasedness, 89–96
- Optimally property, 168
- Orthogonal columns, 326
- Orthogonal polynomials, 311
- Orthogonal projection, 316, 340–341, 509–510
- Orthonormal basis, 404–406
- Orthonormal eigenvectors, 205, 360
- P**
- PACF. *See* Partial autocorrelation function (PACF)
- Parameter estimation
- $ARMA(p, q)$  models, 471–472
  - $AR(p)$  models, 467–469
  - $MA(q)$  models, 469–471
- Parametric and nonparametric models, 70
- Parametric family, 129–131
- Partial autocorrelation function (PACF), 450–453
- ACF, 449–450
- Durbin-Levinson iterative, 451–453
  - Durbin-Levinson recursions, 450–451
- Partial least squares (PLS), 360–361
- Penalty function, 357–358
- Penalty methods, 357–359
- Penalty parameter, 357–358
- Penalty term, 357–358
- Periodogram, 476–478
- Permutation test, 215–216
- Poisson distribution, 206–207, 209
- Polya's theorem, 57–58
- Polynomial model, 311
- Pooled data, 201, 207
- Portmanteau test, 448
- Positively dependent, 151
- Power function, 132
- Prediction, 364
- error, 353–354
  - intervals, 339–340
  - standard error, 465–466
- Principal components analysis, 404–410
- asymptotic results, 408–410
  - estimation of, 407–408
  - orthonormal basis, 404–406
  - regression interpretation, 406–407
- Principal components regression (PCR), 360–361
- Probability analysis
- central limit theorem, 496–500
  - characteristic function, 496–500
  - convex functions, 492–493
  - integration theory, 491–492
  - Stieltjes integral, 493–496
  - weak convergence of, 500–501
  - weak law of large number, 496–500
- Probability, axiomatic definition of, 3
- Probability density function (pdf), 8
- Probability distributions, 69, 169. *See also*
- Continuous distributions, Discrete distributions
  - Probability inequalities, 60–66
  - Probability mass function (pmf), 8
  - Probability space, 3
- Probability theory
- conditional probability and independence, 5–7
  - correlation coefficient, 10–13
  - covariance, 10–13
  - expected value, 10–13
  - moment generating function (mgf), 13
  - moments, 13
  - random experiments, 1
  - set theory, 1–2
  - transforms, 17–21
  - variance, 10–13
- Product-Limit (PL), 273–274
- Projection matrix, 316–317, 353–354
- Proportional reduction, 320, 325–326

- p*-Value  
 accept  $H_0$ /reject  $H_0$ , 159  
 hypothesis testing, 159  
 Pearson's  $P\lambda$  statistic, 159–160  
 test statistic, 159
- Q**  
 Quadratic discriminant rule, 414–417  
 Bayes' rule, 415  
 normal case, 415–416
- Quadratic forms distribution, 510–511
- p*-Quantile, 283
- R**  
 Random-effects model, 361–362  
 Random errors, 310–311  
 Random right-censoring  
   integrated hazard function, 275  
   right algorithm, redistribute to, 274  
   survival function, estimation of, 273–274  
   variance, 274  
 Random stopping time, 166  
 Random variables, transforms of  
   extension, 20  
   linear transformation, 21  
   order statistics, joint distribution of, 20
- Random walk model, 437
- Ranks and order statistics  
   asymptotic distribution, 219–223  
   contiguity theory, 226–227  
   exact distribution under  $H_0$ , 218–219  
   permutation test, 215–216  
   Pitman's approach, asymptotic comparison of, 224–226  
   three basic problems, nonparametric tests in, 216–218
- Rank tests  
   approximate scores, 236–237  
   bivariate population, 234–236  
   general alternative, 228–231  
   LMP, asymptotic distribution of, 237–241  
   one-sample location problem, 232–233  
   two-sample scale problem, 233–234
- Rao-Blackwell formula, 92–93
- Rao-Blackwell method, 94, 115
- Rao-Blackwell theorem, 90
- Rao's statistic, 211
- Real-valued function, 278, 497
- Rectangular density, 186
- Regression  
   partial least squares, 360–361  
   penalty methods, 357–359  
   stepwise, 356–357  
   subsets, 356
- Regression analysis, 406–407
- Regression estimation  
   conditional moments, 272  
   cross-validated choice, 272  
   leave-one-out cross-validation, 272  
   nonparametric estimate, 272  
   optimal choice, 272
- Regression function, 255
- Regression model, 15
- Regression sum of squares, 318–319, 321
- Regular estimators, 192–193
- Regularity conditions, 180–181, 228–229
- Remainder term  $R_n$ , 294–298
- Residual sum of squares, 317–319, 342
- Residual variance, 15
- Restricted maximum likelihood (REML), 372–373
- Ridge regression method, 357–358
- Riemann–Stieltjes integral, 493–496
- Right algorithm, 274
- Right-tail tests, 160
- Risk function, 72
- S**  
 Sample covariance matrix, 387  
 Sample eigenvalues, 408–410  
 Sample mean, 41–42, 170  
 Sample variance, 41–42, 170
- Sampling distributions  
   factorization theorem, 388  
   likelihood function ( $L$ ), 387
- Scheffé method, 389, 395
- Cauchy-Schwarz inequality, 338
- confidence coefficient, 338
- confidence ellipsoid, 337–338
- simultaneous confidence intervals, 337

- Score function  
 asymptotic distribution, 291, 293  
 $L$ -estimator, 286, 291, 293  
 $M$ -estimator, 293  
 $M$ -functional, 296  
 Sequential analysis, 160. *See also* sequential probability ratio test (SPRT)  
 Sequential probability ratio test (SPRT)  
 ASN function, 166–168  
 definition of, 161–162  
 error probabilities of, 163–164  
 OC function, 164–166  
 stops with probability 1, 162–163  
 Set theory, 1–2  
 Sherman-Morrison formula, 354–355, 504  
 Simple hypothesis, 126  
 Simple linear regression model, 157, 309  
 Simple null hypothesis *vs.* simple alternative  
   distinct probability distributions, 127–128  
   existence, 127  
   necessity, 127  
 Neyman-Pearson lemma, 127  
   sufficiency, 127  
 Simple propositions, 3–4  
 Simple *vs.* simple likelihood ratio, 195  
 Simultaneous confidence intervals, 336–337  
 Single-parameter exponential family, 133, 135  
 Singular value decomposition, 505  
 SLLN. *See* Strong law of large numbers (SLLN)  
 Slutsky's theorem, 57–58  
 Smoothing parameter, 255–256  
 Spearman's rank correlation, 218  
 Spectral analysis, 473–486  
   linear filtering, 480–483  
   periodogram, 476–478  
   remarks, 474–475, 479  
   spectral density, 478–480, 483–486  
   stationary series, 475–476  
 Spectral decomposition theorem, 505  
 Spectral density function  
   ARMA, 483–486  
   autocovariance function, 437–440  
   estimation of, 478–480  
   frequency response function, 483–485  
   special cases, 485–486  
 Spectral distribution function, 440  
 SPRT. *See* Sequential probability ratio test (SPRT)  
 Squared-error loss, 89–96  
 Square root of a matrix, 506  
 Standard deviation, 11  
 Standard normal distribution, 36  
 Stationarity  
   ARIMA, 436  
   autocovariance function, 437–440  
   autoregressive process, 436  
   Cramér representation, 475–476  
   Gaussian, 434–435  
   linear prediction, 441–442  
   linear time series, 440–441  
   moving average process, 435  
   random walk model, 437  
   strict, 434  
   time reversibility, 441–442  
   weakly, 435  
   white noise, 435  
 Statistical data analysis, 309  
 Statistical functionals  
   bootstrap method, 298–303  
   delta method, 278–285  
   exercises, 303–307  
   jackknife method, 298–303  
    $L$ -estimators, 285–289  
    $M$ -estimator, 289–292  
   remainder term  $R_n$ , 294–298  
 Statistical inference  
   confidence sets, 71  
   hypothesis testing, 71  
   optimal decision rules, 76–86  
   parametric and nonparametric models, 70  
   point estimation, 71  
   population and random samples, 69  
   problems of, 70–71  
   statistical decision functions, 71–74  
   sufficient statistics, 74–76  
 Stepwise regression  
   backward elimination, 356–357  
   forward selection, 356  
 Stieltjes integrals, properties, 494–496  
 Stirling's approximation, 43

Strictly stationary, 434  
 Strong consistency, 185. *See also* Almost sure convergence  
 Strong law of large numbers (SLLN), 56  
 Strong uniform consistency, 258, 261  
 Studentized version, 170  
 Subsets regression, 356  
 Sufficiency, 179  
 Sufficient statistic, 74–75  
 Sum of squares and products (SSP), 396–397  
 Superefficient estimators, 192  
 Supporting hyperplane theorem, 492  
 Survival function, 273–275  
 Survival time, 273  
 Symmetric matrix function, 506  
     inverse of a matrix, 506  
     inverse of square root of a matrix, 506  
     multivariate analysis, 506  
     square root of a matrix, 506

**T**  
 Taylor series, 278  
 Tchebyshev's inequality, 57, 61  
 $t$  distributions, 42–43  
 Test statistic, 125, 390  
 Time reversibility, 441–442  
 Time series  
     ARMA model appropriate, 472–473  
     autocorrelation function, 442–450  
     autocovariance function, 434, 437–440  
     causality, 454–460  
     forecasting, 460–466  
     invertibility, 454–460  
     mean, 442–450  
     PACE, 450–453  
     parameter estimation, 467–472  
     spectral analysis, 473–486  
     stationarity, 434–442  
 Traditional statistical inference, 160  
 Treatment sum of squares (SSTR), 318–319  
 Triangular density, 186  
 Tukey method  
     application of, 339

studentized range distribution, 339  
 studentized range variable, 339  
 Two-factor ANOVA model, 313–314, 321  
 Two-parameter exponential family, 150  
 Two random factors, 362  
 Two-sided problems, 136–138  
 Two-term Taylor expansion, 266–267  
 Type I error, 125  
 Type I error probability, 78  
 Type II error, 125  
 Type II error probability, 77–78

**U**  
 Unbiased confidence interval, 172–173  
 Unbiased estimators, 182, 191–192  
 Unbiasedness, 89–90  
 Unbiased tests, 78  
     behavioral test, 133  
     MLR property, 133  
     natural parameter space, 133–135  
     null hypothesis, 133  
     power functions, 133  
     single-parameter exponential family, 133  
 Uniform distribution, 34  
 Uniform integrability, 56, 350–351  
 Uniformly minimum variance unbiased estimator (UMVUE), 89–96  
 Uniformly most powerful (UMP)  
     nuisance parameters, 144–158  
     one-sided hypotheses, 131–132  
     two-sided problems, unbiased tests, 136–139  
 Upper confidence bounds, 169  
 $U$ -statistic, 277

**V**  
 Variance components  
     Henderson's method III, 369–371  
     maximum likelihood, 369, 371–373  
     MINQUE, 369  
     restricted maximum likelihood, 369  
 Variance-stabilizing transformations, 60  
 $V$ -statistic, 277, 284

**W**

Wald's statistics, 211

Weak convergence, probabilities,  
500–501

Weak law of large numbers (WLLN), 56

Wilcoxon signed-rank statistic, 216–217

Wilks' lambda, 396

exact distributions, 397–398

factor B main effects, 401

Wishart distribution, 383–385

WLLN. *See* Weak law of large numbers (WLLN)

**Y**

Yule-Walker equations, 441–442, 452–453