

Studies in Theoretical and Applied Statistics
Selected Papers of the Statistical Societies

Giorgio Alleva
Andrea Giommi *Editors*

Topics in Theoretical and Applied Statistics

 Springer

Studies in Theoretical and Applied Statistics

Selected Papers of the Statistical Societies

Series Editors

Societa Italiana di Statistica (SIS)

Spanish Society of Statistics and Operations Research (SEIO)

Société Française de Statistique (SFdS)

Sociedade Portuguesa de Estatística (SPE)

Federation of European National Statistical Societies (FENStatS)

More information about this series at <http://www.springer.com/series/10104>

Giorgio Alleva • Andrea Giommi
Editors

Topics in Theoretical and Applied Statistics

 Springer

Editors

Giorgio Alleva
MEMOTEF
Sapienza University of Rome
Rome, Italy

Andrea Giommi
Dept. of Statistics & Informatics
University of Florence
Florence, Italy

ISSN 2194-7767 ISSN 2194-7775 (electronic)
Studies in Theoretical and Applied Statistics
ISBN 978-3-319-27272-6 ISBN 978-3-319-27274-0 (eBook)
DOI 10.1007/978-3-319-27274-0

Library of Congress Control Number: 2016932383

Springer Cham Heidelberg New York Dordrecht London
© Springer International Publishing Switzerland 2016

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

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

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media
(www.springer.com)

Foreword

Dear reader,

On behalf of the four Scientific Statistical Societies—the SEIO, Sociedad de Estadística e Investigación Operativa (Spanish Society of Statistics and Operations Research); SFdS, Société Française de Statistique (French Statistical Society); SIS, Società Italiana di Statistica (Italian Statistical Society); and the SPE, Sociedade Portuguesa de Estatística (Portuguese Statistical Society)—we would like to inform you that this is a new book series of Springer entitled *Studies in Theoretical and Applied Statistics*, with two lines of books published in the series: *Advanced Studies* and *Selected Papers of the Statistical Societies*.

The first line of books offers constant up-to-date information on the most recent developments and methods in the fields of theoretical statistics, applied statistics, and demography. Books in this series are solicited in constant cooperation between the statistical societies and need to show a high-level authorship formed by a team preferably from different groups so as to integrate different research perspectives.

The second line of books presents a fully peer-reviewed selection of papers on specific relevant topics organized by the editors, also on the occasion of conferences, to show their research directions and developments in important topics, quickly and informally, but with a high level of quality. The explicit aim is to summarize and communicate current knowledge in an accessible way. This line of books will not include conference proceedings and will strive to become a premier communication medium in the scientific statistical community by receiving an impact factor, as have other book series such as *Lecture Notes in Mathematics*. The volumes of selected papers from the statistical societies will cover a broad range of theoretical, methodological, as well as application-oriented articles, surveys, and discussions. A major goal is to show the intensive interplay between various, seemingly unrelated domains and to foster the cooperation between scientists in different fields by offering well-founded and innovative solutions to urgent practice-related problems.

On behalf of the founding statistical societies, I wish to thank Springer, Heidelberg, and in particular Dr. Martina Bihn for the help and constant cooperation in the organization of this new and innovative book series.

Rome, Italy

Maurizio Vichi

Preface

The Italian Statistical Society (SIS) holds a biyearly international scientific meeting where both methodological and applied statistical research papers are welcome. The aim of this volume is to present significant and innovative contributions which were presented in a preliminary version at the 46th International Meeting of the SIS in Rome. More than 250 contributions were presented at this meeting by about 500 scientists and experts coming from several countries. Fifty-eight extended versions of these contributions were subsequently submitted for potential inclusion in this volume. After a careful double-blind review process, carried out with the help of approximately one hundred referees, 27 of these papers were chosen.

The volume is organized into five parts: the first three are in prevalence based on a methodological framework; and the other two regard applied issues. The first part collects miscellaneous contributions on statistical theory. The second focuses on methods for data mining and multivariate data analysis. The papers included in the third part deal with sampling and estimation methods. The papers in the fourth part focus on the application of statistical methods in the analysis of social, demographic, and health data, and the last part is dedicated to the analysis of economic and econometric features.

The editors are grateful to all the referees for their conscientious work. Finally, special thanks go to Alice Blanck from Springer Verlag for her patience and valued assistance in preparing this volume.

Rome, Italy
Florence, Italy

Giorgio Alleva
Andrea Giommi

Contents

Part I Statistical Theory and Methods

Empirical Orthogonal Function and Functional Data Analysis Procedures to Impute Long Gaps in Environmental Data	3
Francesca Di Salvo, Antonella Plaia, Mariantonietta Ruggieri, and Gianna Agró	
Unconditional and Conditional Quantile Treatment Effect: Identification Strategies and Interpretations	15
Margherita Fort	
Some New Results on the Beta Skew-Normal Distribution	25
Valentina Mamelì and Monica Musio	
The Median of a Set of Histogram Data	37
Lidia Rivoli, Rosanna Verde, and Antonio Irpino	
Rates for Bayesian Estimation of Location-Scale Mixtures of Super-Smooth Densities	49
Catia Scricciolo	

Part II Data Mining and Multivariate Data Analysis

Unsupervised Classification of Multivariate Time Series Data for the Identification of Sea Regimes	61
Mauro Bencivenga, Francesco Lagona, Antonello Maruotti, Gabriele Nardone, and Marco Picone	
An Evaluation of the Student Satisfaction Based on CUB Models	73
Barbara Cafarelli and Corrado Crocetta	
Dimensions of Well-Being and Their Statistical Measurements	85
Carla Ferrara, Francesca Martella, and Maurizio Vichi	
Extracting Meta-information by Using Network Analysis Tools	101
Agnieszka Stawinoga, Maria Spano, and Nicole Triunfo	

Factor PD-Co-clustering on Textual Data	111
Cristina Tortora, Marina Marino, and Germana Scepi	
Part III Sampling and Estimation Methods	
M-Quantile Small Area Estimation for Panel Data	123
Annamaria Bianchi	
A Two-Part Geoadditive Small Area Model for Geographical Domain Estimation	133
Chiara Bocci, Alessandra Petrucci, and Emilia Rocco	
A Unified Approach for Defining Optimal Multivariate and Multi-Domains Sampling Designs	145
Piero Demetrio Falorsi and Paolo Righi	
Estimation of Poverty Rate and Quintile Share Ratio for Domains and Small Areas	153
Risto Lehtonen and Ari Veijanen	
A Sample Survey on Inactive Students: Weighting Issues in Modelling the Inactivity Status	167
Lucio Masserini and Monica Pratesi	
Part IV Social Statistics, Demography and Health Data	
The Material Deprivation of Foreigners: Measurement and Determinants	181
Annalisa Busetta, Anna Maria Milito, and Antonino Mario Oliveri	
How Do Life Course Events Affect Paid and Unpaid Work of Italian Couples?	193
Maria Gabriella Campolo, Antonino Di Pino, and Ester Lucia Rizzi	
Do Rational Choices Guide Family Formation and Dissolution in Italy?	205
Gustavo De Santis and Silvana Salvini	
STAR Modeling of Pulmonary Tuberculosis Delay-Time in Diagnosis	215
Bruno de Sousa, Dulce Gomes, Patrícia A. Filipe, Cristiana Areias, Teodoro Briz, Carlos Pires, and Carla Nunes	
Non-aggregative Assessment of Subjective Well-Being	227
Marco Fattore, Filomena Maggino, and Alberto Arcagni	
Composite Indicator of Social Inclusion for the EU Countries	239
Francesca Giambona and Erasmo Vassallo	

A Well-Being Index Based on the Weighted Product Method	253
Matteo Mazziotta and Adriano Pareto	
 Part V Economic Statistics and Econometrics	
 A Comparison of Different Procedures for Combining High-Dimensional Multivariate Volatility Forecasts	263
Alessandra Amendola and Giuseppe Storti	
 Which Seasonality in Italian Daily Electricity Prices? A Study with State Space Models	275
Paolo Chirico	
 From the Standard of Living as a Latent Variable to the Estimation of Equivalence Scales and Other Indices	285
Gustavo De Santis and Mauro Maltagliati	
 Learning-by-Exporting and Productivity: Evidences from a Panel of Manufacturing Firms	295
Maria Rosaria Ferrante, Marzia Freo, and Alessandro Viviani	
 A Multivariate VEC-BEKK Model for Portfolio Selection	307
Andrea Pierini and Alessia Naccarato	

Part I

Statistical Theory and Methods

Empirical Orthogonal Function and Functional Data Analysis Procedures to Impute Long Gaps in Environmental Data

Francesca Di Salvo, Antonella Plaia, Mariantonietta Ruggieri,
and Gianna Agró

Abstract

Air pollution data sets are usually spatio-temporal multivariate data related to time series of different pollutants recorded by a monitoring network.

To improve the estimate of functional data when missing values, and mainly long gaps, are present in the original data set, some procedures are here proposed considering jointly Functional Data Analysis and Empirical Orthogonal Function approaches. In order to compare and validate the proposed procedures, a simulation plan is carried out and some performance indicators are computed. The obtained results show that one of the proposed procedures works better than the others, providing a better reconstruction especially in presence of long gaps.

1 Introduction

Imputing missing values is a very crucial issue in many fields [4], especially in air pollution data sets, where often high percentages of data are missing and long gap sequences may occur, due to failures of monitoring instruments or integration of data from mobile and fixed stations. At this aim many methods have been proposed in literature, such as *Kriging* and other *optimal interpolation* methods, for example, *objective analysis*. Both Empirical Orthogonal Function (EOF) and Functional Data Analysis (FDA) are also used for imputing missing values and as denoising tools at the same time. In particular, EOF methodology is one of the emerging approaches

F. Di Salvo • A. Plaia (✉) • M. Ruggieri • G. Agró
Department of Statistical and Mathematical Sciences, University of Palermo, Viale delle Scienze,
Ed. 13, 90128 Palermo, Italy
e-mail: francesca.disalvo@unipa.it; antonella.plaia@unipa.it; mariantonietta.ruggieri@unipa.it;
gianna.agro@unipa.it

in this framework; it has been widely used for oceanographic and meteorological applications to fill in missing data in spatio-temporal *univariate* data sets and it is particularly valid when a high percentage of data is missing [1, 12].

The objective of this paper is to extend the EOF to a spatio-temporal *multivariate* data set, described as *functional*, to improve the estimate of functional with a reconstruction of signal in correspondence of long gaps. Two procedures, performing EOF and FDA jointly, are here proposed and compared:

- P1: EOFs are computed on the observed data by a Singular Value Decomposition (SVD) and a Principal Component Analysis (PCA) and then the approximated data are converted into functional;
- P2: data are converted into functional and the reconstruction is obtained by Functional Singular Value Decomposition (FSVD) and Functional Principal Component Analysis (FPCA).

The two procedures are compared by computing two performance indicators [10] on simulated missing data. The entire analysis is implemented in R by using also *fda* package (<http://cran.r-project.org>).

In Sect. 2 the observed data set is described, while in Sect. 3 the main characteristics of the FDA and the EOF approaches are briefly outlined. In Sect. 4 the simulation plan and the initial missing values imputation are presented, then the proposed procedures are introduced. Finally, in Sect. 5 the obtained results are shown and in Sect. 6 some conclusions are drawn.

2 The Air Pollution Data Set

A spatio-temporal *multivariate* data set, related to four main pollutant concentrations (CO, NO₂, PM₁₀ and SO₂) hourly (or bi-hourly) recorded at the nine monitoring stations of Palermo (Italy) during 2005, is considered. O₃ is not included in the analysis, being available only for two stations. Data are provided by AMIA (Azienda Municipalizzata Igiene Ambientale, <http://www.amianet.it/>).

To obtain daily syntheses, data are aggregated by time, at each site for each pollutant, using the functions suggested by EC guidelines [2].

To compare pollutants, different for measurement unit or order of magnitude, data are also standardized. Among all the standardizing transformations proposed in literature, we prefer the linear interpolation [6] and use the EC directive thresholds reported in [5], obtaining values in [0, 100]. Details about the reason of this choice are reported in [9].

Data are organized in a three way ($T \times N \times P = 365 \times 9 \times 4$) array of standardized observed data.

3 The Approaches Based on FDA and EOF

3.1 The Empirical Orthogonal Functions

The EOF analysis has become an interesting statistical tool for reconstructing the data, through the identification of structures in space–time relationships [1, 12].

Performing an SVD, the focus is extracting empirically salient modes of variation from the spatial and temporal singular vectors of the data matrix. Starting from a three-way array $T \times N \times P$, we arrange the observation of a pollutant p_j referred to N sites and T instances of time, into the matrix $\mathbf{X}_{(N \times T)}^{p_j}$.

Since we aim to account for joint multivariate spatio-temporal covariance structures, in our approach the SVD is performed on the matrix

$$\mathbf{X} = [\mathbf{X}^{p_1}, \mathbf{X}^{p_2}, \dots, \mathbf{X}^{p_P}].$$

The i th row of X gives the profile of the data in the i th site, with the T columns of the p th submatrix relating to timing variations in the p th pollutant. The generic t th column of the p th submatrix gives the space profile of the p th pollutant in a single day.

SVD decomposes the spatio-temporal data matrix in the product $\mathbf{X} = \mathbf{U}\mathbf{\Gamma}\mathbf{A}$, where $\mathbf{U} = \{u_{ir}\}$ is the matrix of left singular vectors, spanning the space of the variables along the sites, $\mathbf{A} = \{a_{rt}^{pj}\}$ is the matrix of the right singular vectors spanning the space of the variables along the time, and $\mathbf{\Gamma}$ is the diagonal matrix with elements $\{\sqrt{\gamma_r}\}$.

By selecting the first $v \leq \text{rank}(\mathbf{X})$ eigenvalues of the matrix $\mathbf{\Gamma}$:

$$\mathbf{X}_{N \times (T \times P)} \approx \mathbf{U}_{N \times v} \mathbf{\Gamma}_{v \times v} \mathbf{A}_{v \times (T \times P)},$$

the data are given the interpretation of the sum of multiplicative effects of spatial and temporal factors:

$$x_i^{pj}(t) \approx \sum_{r=1}^v \gamma_r u_{ir} a_{rt}^{pj}. \quad (1)$$

The relationship with PCA is well known when the principal components are computed from the covariance matrix of \mathbf{X} .

Since the pollutants are correlated within each station [11], the approach allows to exploit and recover their simultaneous variability.

3.2 Functional EOFs

The advantage of using FDA depends on the nature of observed data; such an approach is just suggested by the functional structure of the data. In FDA observed data are considered as continuous functions of time sampled in discrete times [7] and subject to observational noise. The error term is subject to the usual iid assumptions with zero mean and constant finite variance; converting observed data into functional aims also at removing such a random component.

In this context, the generic realization x_{is}^{Pj} , recorded at time s ($s = 1, \dots, T$) for the pollutant p_j ($j = 1, \dots, P$) at the station i ($i = 1, \dots, N$), is the result of a signal, $\tilde{x}_i^{Pj}(t)$, affected by a noise ε_{is}^{Pj} :

$$x_{is}^{Pj} = \tilde{x}_i^{Pj}(t) + \varepsilon_{is}^{Pj}. \quad (2)$$

The curves $\tilde{x}_i^{Pj}(t)$ in (2) may be expressed in terms of a linear combination of a complete set of K suitable basis functions ϕ_k ; details about smoothing strategies for functional data are exhaustively treated in [8]. Here we describe the adopted strategy: a unique basis system is chosen for all the pollutants, since in a multivariate approach, we are taking into account the simultaneous variability of all the pollutants:

$$\tilde{x}_i^{Pj}(t) = \sum_k^K c_{i,k}^{Pj} \phi_k(t). \quad (3)$$

Making use of regression splines, which are functions obtained by joining segments of polynomials smoothed at points called knots, we avoid to impose uniform cyclicity on the curves. The implementation makes use of cubic B-spline basis system with $K = 179$ equally spaced knots, that allow to capture seasonal, monthly and weekly variations [3], but also events that occur irregularly and cannot be expected to be periodically repeated.

In the (3) the $N \times K \times P$ coefficients $c_{i,k}^{Pj}$ are estimated by penalized least squares (see [8]). Here the chosen value for the smoothing parameter ($\lambda = 2$) appears to be a fair compromise between what can be suggested by an automatic method, such as the generalized cross validation, and a subjective choice, that aims at smoothing rough data without hiding their variability linked to possible peaks. On the whole, our choices seem to be a good trade-off in smoothing between the removal of measurement error and the preservation of information.

For the second Procedure P2, EOF is carried out on the functional data set $\tilde{\mathbf{X}}$, with the advantages of dealing with a few coefficients, rather than a large number of data [8].

The FPCA, which takes place in the space spanned by the basis function set, was deeply studied by Ramsay and Silverman [8]; in [11] we discussed the computational aspects for the extension to the multipollutant case and we also developed the computational steps for the FSVD.

Due to the high correlations among the curves of the pollutants in each site and the high correlations among the sites for a single pollutant, our approach allows to capture the temporal dynamic of a pollutant in the whole area and the local dynamics in a site for all the pollutants. The crucial node of this approach is that the generalization to the functional multivariate setting must preserve the well known PCA and SVD properties: the appeal of an FSVD must consist in the ability to extract the most relevant mode of variations from the spatial and temporal singular vectors of the matrix whose elements are defined in terms of the basis functions and coefficients.

As detailed in [11] we move from the linear expansion (3) of the curves, defining the matrices \mathbf{Z}^{pj} for each pollutant:

$$\mathbf{Z}^{pj} = N^{-1} \mathbf{C}^{pj} \mathbf{W} \mathbf{L},$$

where $\mathbf{C}_{(N \times K)}^{pj}$ is the matrix of coefficients $c_{i,k}^{pj}$, \mathbf{W} is the order K symmetric matrix in terms of the basis system Φ_k :

$$\mathbf{W} = \int \Phi(t) \Phi(t)' dt, \quad (4)$$

and \mathbf{L} is the inverse matrix of the Cholesky decomposition of \mathbf{W} , i.e. $\mathbf{L} = (\mathbf{W}^{\frac{1}{2}})^{-1}$. Accounting for joint multivariate spatio-temporal covariance structures, in our approach the FSVD is based on the SVD of the matrix \mathbf{Z} :

$$\mathbf{Z} = [\mathbf{Z}^{p1}, \mathbf{Z}^{p2}, \dots, \mathbf{Z}^{pP}].$$

The i th row of \mathbf{Z} gives the temporal profile of P curves in the i th site. The k th column for the p th submatrix gives the spatial profile of the coefficients along the N sites for the k th basis and the p th pollutant.

FSVD decomposes \mathbf{Z} in the product:

$$\mathbf{Z}_{N \times (K \times P)} = \mathbf{U}_{N \times (K \times P)} \Gamma_{(K \times P) \times (K \times P)} \mathbf{A}_{(K \times P) \times (K \times P)},$$

where the columns of \mathbf{U} , $\{\mathbf{u}_r\}$ are the left singular vectors, spanning the space of the coefficients along the sites, the rows of \mathbf{A} , $\{\mathbf{a}_r\}$, are the right singular vectors spanning the space of the variables along the time, and Γ is the diagonal matrix with elements $\{\sqrt{\gamma_r}\}$.

By selecting the first $\nu \leq \text{rank}(\mathbf{Z})$ eigenvalues of the matrix Γ :

$$\mathbf{Z}_{N \times (K \times P)} \approx \mathbf{U}_{N \times \nu} \Gamma_{\nu \times \nu} \mathbf{A}_{\nu \times (K \times P)} \quad (5)$$

the best rank- ν approximated coefficients can be interpreted in terms of the multiplicative effects of spatial and temporal factors. Splitting the matrix $\mathbf{A}_{\nu \times (K \times P)}$

into P parts:

$$\mathbf{A}_{v \times (K \times P)} = \left[\mathbf{A}_{v \times K}^{p1}, \dots, \mathbf{A}_{v \times K}^{pP} \right]$$

we obtain:

$$z_{ik}^{pj} \approx \sum_{r=1}^v \gamma_r u_{ir} a_{(rk)}^{pj}. \quad (6)$$

The approximated matrix, obtained by (6), contains new coefficients for the basis system in the linear expansion of the data.

The FPCA also decomposes variability in functional data, finding directions along which data have the highest variability. This goal is achieved by an eigenanalysis of the variance operator yielding eigenfunctions $\xi_m(t)$ that vary with time, after having defined approximations of the continuous eigenfunctions in terms of the centred smoothed functional data. The eigenfunctions form a set of EOFs (in [8] they are named *harmonics*). The principal scores, uncorrelated in the new coordinate system, are weighted at each instance of time by the EOFs, determining the approximated rank- v matrix, whose elements are the coefficients for the basis system in the linear expansion of the data; from the relationship between FPCA and FSVD, it follows that the standardized principal scores are the columns of the matrix \mathbf{U} (see [11] for a detailed description).

A potential weakness of this reconstruction is that it is optimal only when the underlying pattern is linear, while nonlinear processes in general may be more realistic; however, it provides a good approximation and, because of the orthogonality of the directions of variation, simple interpretations.

4 The Initial Imputation and the Proposed Procedures

In order to compare the performance of different imputation methods, many studies use real data sets and simulated missing data patterns by deleting values.

In this paper, to validate the proposed imputation procedures, simulated incomplete data sets are generated reproducing the actual pattern of missing data, that is the same pattern of the observed data set. In particular, 100 missing data indicator arrays $\mathbf{M}_{T \times N \times P}$ ($365 \times 9 \times 4$), with dimensions as our data set, are randomly generated from a Bernoulli distribution with parameter π equal to the actual percentage of missing in each monitoring station. It can often happen that very long gap sequences are observed in air pollution data sets, due to long time failures not easily solvable or to data coming from a mobile monitoring station, therefore some gap sequences of 2 or 3 months long are also randomly generated and randomly placed, according to pollutant and station, in each array \mathbf{M} . Then, each array \mathbf{M} is applied to the observed data set $\mathbf{X}_{365 \times 9 \times 4}$, creating “artificial” gaps and obtaining 100 arrays \mathbf{X}^M . Since the values corresponding to the gaps artificially created are known, the computation

of some performance indicators allows to assess the goodness of the imputation methods.

Before performing any procedure on each \mathbf{X}^M , missing data must be replaced by some initial values: fixing the pollutant, they are first filled by a *station mean* (Ms), that is the annual mean for each station or by a *day mean* (Md), that is the mean among stations for each day.

After filling initially artificial missing by Ms or Md , the two proposed procedures are carried out on each \mathbf{X}^M ; in order not to lose any information, only missing values are replaced.

In particular:

- **Procedure 1 (P1)**

The EOF (PCA and SVD) is performed on each \mathbf{X}^M and the reconstructed array is converted into functional obtaining $\tilde{\mathbf{X}}_{P1}^M$;

- **Procedure 2 (P2)**

Each array \mathbf{X}^M is converted into the functional $\tilde{\mathbf{X}}^M$ and the functional EOF (FPCA and FSVD) is performed on $\tilde{\mathbf{X}}^M$ obtaining $\tilde{\mathbf{X}}_{P2}^M$.

The aim is to investigate if a preliminary functional transformation allows a better reconstruction by EOF in presence of long gap sequences.

The optimal number ν of EOFs to be extracted can be determined in different ways, for example, by cross-validation techniques. In this paper the proposed procedures are compared with the same ν ; such a value is chosen on the basis of the explained variability (more than 95 %).

After obtaining the reconstructed arrays $\tilde{\mathbf{X}}_{P1}^M$ and $\tilde{\mathbf{X}}_{P2}^M$ by the procedures P1 and P2, these are compared with $\tilde{\mathbf{X}}$, representing actual data, by means of two performance indicators (for a detailed description see [10]) :

- the correlation coefficient ρ ,
- the root mean square deviation $RMSD$.

ρ and $RMSD$ are computed between the set of missing values or just long gap sequences (and not the whole arrays). Of course, the higher is ρ and the lower is $RMSD$, the better is the reconstruction. $RMSD$, with respect to ρ , is related to the size of the discrepancies between actual and imputed values, while ρ is related to data variability.

5 Results

The distributions of the two performance indicators ρ and $RMSD$, over the 100 arrays \mathbf{X}^M , are summarized by their means M and standard deviations σ considering, for the two different initial imputations Md and Ms , the whole set of missing values or only long gap sequences (Table 1). Actually, due to seasonal behaviour of a

Table 1 M and σ of performance indicators

	Initial imputation: Md				Initial imputation: Ms			
	$P1_{PCA}$	$P1_{SVD}$	$P2_{FPCA}$	$P2_{SVD}$	$P1_{PCA}$	$P1_{SVD}$	$P2_{FPCA}$	$P2_{SVD}$
<i>All missing values</i>								
M_ρ	0.9724	0.9742	0.9708	0.9568	0.9773	0.9795	0.9756	0.9690
σ_ρ	0.0122	0.0125	0.0121	0.0126	0.0070	0.0070	0.0069	0.0070
M_{RMSD}	3.3194	3.1928	3.4149	4.1434	2.9852	2.8312	3.0969	3.5357
σ_{RMSD}	0.7136	0.7445	0.6891	0.5841	0.4172	0.4279	0.3903	0.3147
<i>Only long gap sequences</i>								
M_ρ	0.7747	0.7716	0.7722	0.7779	0.5137	0.5876	0.5238	0.7478
σ_ρ	0.2341	0.2399	0.2365	0.2365	0.3272	0.3446	0.3155	0.2248
M_{RMSD}	3.9822	3.8306	4.0046	4.1660	3.8386	3.5845	3.8574	3.5852
σ_{RMSD}	2.5647	2.5391	2.5507	2.6019	2.1292	2.0711	2.1069	1.7818

pollutant, Ms could be a bad initial imputation, nevertheless it is here considered to test the robustness of the proposed procedures.

As it can be observed, when the initial imputation is the day mean Md , the results obtained on all missing values by the four procedures are quite similar for both performance indicators; the same happens if only long gap sequences are considered.

When the initial imputation is the station mean Ms , the results obtained on all missing values by the four procedures are still quite similar for both performance indicators but, if only long gap sequences are considered, $P2_{SVD}$ outperforms the other three procedures (Table 1 in bold), giving results comparable to the ones got when the initial imputation is the day mean Md . Thus, $P2_{SVD}$ should seem not to be influenced by the initial imputation, appearing more “robust” with respect to the other procedures. Actually, $P2_{SVD}$ and $P1_{SVD}$ provide in average the same results for $RMSD$, but the standard deviation of $RMSD$ is lower for $P2_{SVD}$.

What we now claim is also evident by the distributions of the two performance indicators over the 100 arrays \mathbf{X}^M , reported in Fig. 1 and obtained considering Ms as initial imputation and only long gap sequences. The ρ distribution is always negatively asymmetric, but it shows a noticeable higher percentage of arrays with a higher ρ for $P2_{SVD}$. The $RMSD$ distribution appears positively asymmetric for all the four procedures, confirming their good performance, but presents a lower variability for $P2_{SVD}$.

In Fig. 2 some examples of long gaps reconstructed by the four procedures are shown. As it can be noted, the curves obtained by $P2_{SVD}$ are very close to the curves FD , representing actual observed data after denoising; these lines seem to follow the continuous ones by catching the same peaks. So, the procedure $P2_{SVD}$ enables to recover the data variability by taking into account the recorded peaks better than

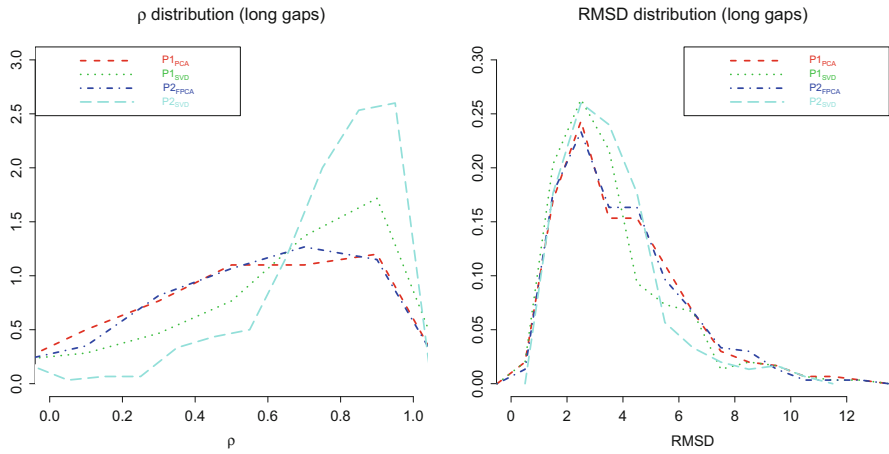


Fig. 1 Performance indicator distributions (initial imputation: M_S)

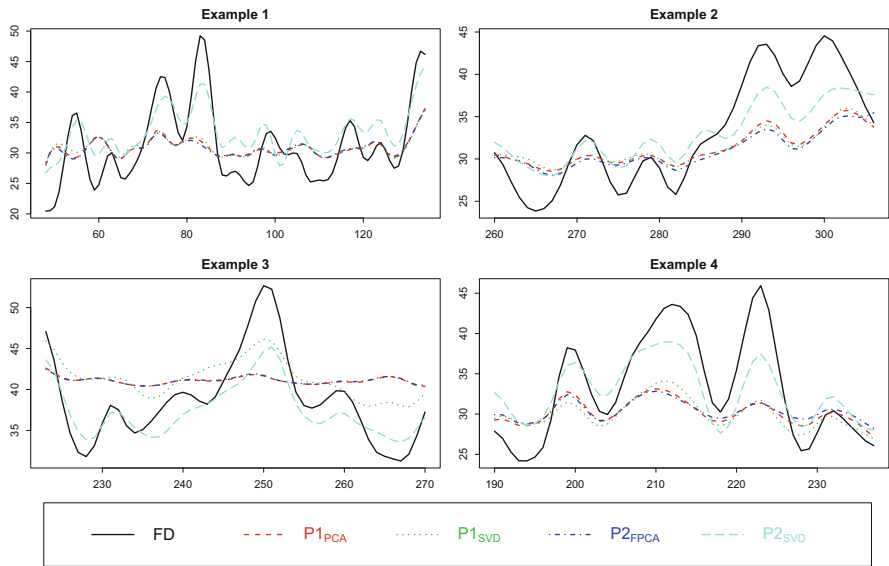


Fig. 2 Examples of long gaps reconstructed by the proposed procedures (initial imputation: M_S)

the other procedures. This is also confirmed by the values of $RMSD$ and ρ , reported in Table 2, lower and higher for $P2_{SVD}$, respectively.

Table 2 Performance indicators for the examples reported in Fig. 2

	ρ				RMSD			
	$P1_{PCA}$	$P1_{SVD}$	$P2_{FPCA}$	$P2_{SVD}$	$P1_{PCA}$	$P1_{SVD}$	$P2_{FPCA}$	$P2_{SVD}$
Example 1	0.6335	0.6512	0.5746	0.9261	5.6344	5.6076	5.7629	3.9345
Example 2	0.8965	0.8502	0.8400	0.9348	5.0189	5.1765	5.2945	3.5141
Example 3	0.4139	0.8462	0.4322	0.9428	5.6838	4.4335	5.6928	2.9402
Example 4	0.8501	0.8533	0.8374	0.8939	5.9912	5.8935	6.1259	3.6078

6 Conclusions

In this work some procedures are compared using a spatio-temporal multivariate data set related to air pollution, in order to find the most effective one for estimating functional data in presence of missing values; in particular, the paper focuses on long gap sequences. To deal with functional rather than raw data it aims also at denoising pollution time series by fluctuations due to contingent factors.

The proposed procedures take jointly into account FDA and EOF. In the Procedure 1 observed data are first reconstructed by EOF and then converted into functional; in the Procedure 2 observed data are first transformed into functional and then EOF reconstruction is applied. The aim is to investigate if a better reconstruction, especially in presence of long gaps, can be provided by EOF after a preliminary pre-processing of data by FDA. Cubic B-spline basis system has been used, nevertheless the procedures result to be robust to the basis choice, since similar results have been observed considering Fourier basis system.

Data sets with simulated pattern of missing are used in order to test the validity of the proposed procedures by means of some performance indicators.

On the basis of the obtained results, the Procedure 2 by FSVD (SVD applied to functional converted data) outperforms the Procedure 1: it seems not to be influenced by the different initial imputations, appearing more “robust” with respect to the Procedure 1 and providing very interesting results in presence of long gap sequences also if an initially raw imputation, as a station mean, is used to replace initially missing values.

Such a feature results to be particularly attractive when long sequences of missing data, occurring from mobile and fixed sites, need to be integrated.

Acknowledgements We would like to thank the referee for his useful suggestions, that we tried to use to improve the paper.

References

1. Beckers, J.M., Rixen, M.: EOF calculations and data filling from incomplete oceanographic datasets. *J. Atmos. Ocean. Technol.* **20**(12), 1839–1856 (2003)
2. European Community: Directive 2008/50/EC of the European Parliament and of the Council of 21 May 2008 on ambient air quality and cleaner air for Europe. Official Journal L 152, 11/6/2008, pp 1–44 (2008)
3. Ignaccolo, R., Ghigo, S., Giovenali, E.: Analysis of air quality monitoring networks by functional clustering. *Environmetrics* **19**, 672–686 (2008)
4. Little, R.J.A., Rubin, D.B.: *Statistical Analysis with Missing Data*. Wiley, New York (1987)
5. Murena, F.: Measuring air quality over large urban areas: development and application of an air pollution index at the urban area of Naples. *Atmos. Environ.* **38**, 6195–6202 (2004)
6. Ott, W.R., Hunt Jr., W.F.: A quantitative evaluation of the pollutant standards index. *J. Air Pollut. Control Assoc.* **26**, 1051–1054 (1976)
7. Ramsay, J.O., Silverman, B.W.: *Applied Functional Data Analysis*. Springer, New York (2002)
8. Ramsay, J.O., Silverman, B.W.: *Functional Data Analysis*, 2nd edn. Springer, New York (2005)
9. Ruggieri, M., Plaia, A.: An aggregate AQI: comparing different standardizations and introducing a variability index. *Sci Total Environ.* (2012). doi:[10.1016/j.scitotenv.2011.09.019](https://doi.org/10.1016/j.scitotenv.2011.09.019)
10. Ruggieri, M., Di Salvo, F., Plaia, A., Agró, G.: EOFs for gap filling in multivariate air quality data: a FDA approach. In: 19th International Conference on Computational Statistics, Paris, 22–27 Aug 2010
11. Ruggieri, M., Plaia, A., Agró, G., Di Salvo, F.: Functional principal component analysis for the explorative analysis of multisite-multivariate air pollution time series with long gaps. *J. Appl. Stat.* **40**(4), 795–807 (2013). doi:[10.1080/02664763.2012.754852](https://doi.org/10.1080/02664763.2012.754852)
12. Sorjamaa, A., Lendasse, A., Cornet, Y., Deleersnijder, E.: An improved methodology for filling missing values in spatiotemporal climate data set. *Comput. Geosci.* **14**(1), 55–64 (2009)

Unconditional and Conditional Quantile Treatment Effect: Identification Strategies and Interpretations

Margherita Fort

Abstract

This paper reviews strategies that allow one to identify the effects of policy interventions on the unconditional or conditional distribution of the outcome of interest. This distinction is irrelevant when one focuses on average treatment effects since identifying assumptions typically do not affect the parameter's interpretation. Conversely, finding the appropriate answer to a research question on the effects over the distribution requires particular attention in the choice of the identification strategy. Indeed, quantiles of the conditional and unconditional distribution of a random variable carry a different meaning even if identification of both these set of parameters may require conditioning on observed covariates.

1 Introduction

In the recent years there has been a growing interest in the evaluation literature for models that allow essential heterogeneity in the treatment parameters and more generally for models that are informative on the impact distribution. The recent increase in the attention devoted to the identification and estimation of quantile treatment effects (QTEs) is due to their intrinsic ability to characterize the heterogeneous impact of the treatment on various points of the outcome distribution. QTEs are informative about the impact distribution when the potential outcomes observed under various levels of the treatment are comonotonic random variables. The variable describing the relative position of an individual in the outcome

M. Fort (✉)

Department of Economics, University of Bologna, IZA and CHILD, Piazza Scaravilli 2, Bologna, Italy

e-mail: margherita.fort@unibo.it

distribution thus plays a special role in this setting, representing at the same time the main dimension along which treatment effects are allowed to vary as well as a key ingredient to relate potential outcomes. Several identification approaches currently used in the literature for the assessment of mean effects have thus been extended to quantiles. Most of these strategies require to condition on a set of variables to achieve identification. While conditioning on a set of observed regressors does not affect the interpretation of the parameters in a mean regression, this is not the case for quantiles. The law of iterated expectations guarantees that the parameters of a mean regression have both a conditional and an unconditional mean interpretation. This does not carry over to quantiles, where conditioning on covariates affects the interpretation of the residual disturbance term. Indeed, since quantile regression allows one to characterize the heterogeneity of the treatment response only along this latter dimension, conditioning on covariates in quantile regression generally affects the interpretation of the results.

This paper reviews strategies aimed at identifying QTEs, covering strategies that deal with the identification of conditional and unconditional QTEs with particular attention to cross-sectional data applications in which the treatment is endogenous without conditioning on additional covariates. The aim of the paper is to provide useful guidance for users of quantile regression methods in choosing the most appropriate approach while addressing a specific research question.

The remainder of the paper is organized as follows. After introducing the basic notation and the key parameters of interest in Sects. 2 and 3 reviews solutions to the identification of QTEs. The review covers strategies that are appropriate only when the outcome of interest is a continuous variable, i.e. in cases where the quantiles of the outcome distribution are unambiguously defined. It concludes illustrating some of the methods through two illustrative examples aimed at assessing the distributional impacts of training on earnings and of education on wages. Section 4 concludes.

2 What Are We After: Notation and Parameters of Interest

In this section I first introduce the notation used throughout the paper and then define the objects whose identification is sought.

Y denotes the observed outcome, D the intensity of the treatment received and W a set of observable individual characteristics. W may include exogenous variables X and instruments Z .¹ Y is restricted to be continuous while D , W can be either continuous or discrete random variables. Both Y and D can be decomposed in two components: one of which is deterministic and one of which is stochastic. These two components need not be additively separable. The stochastic components account for differences in the distribution of D and Y across otherwise identical individuals. The econometric models reviewed in Sect. 3 place restriction on: (1) the scale of

¹Capital letters denote random variables and lower case letters denote realizations.

D ; (2) the number of independent sources of stochastic variation in the model; (3) the distribution (joint, marginal, conditional) of these stochastic components and D or $W \equiv (X, Z)$; (4) the scale of Z . Y_i^d denotes the *potential outcome* for individual i if the value of the treatment is d : it represents the outcome that would be observed had the individual i been exposed to level d of the treatment. $F_{Y^d}(\cdot)$, $f_{Y^d}(\cdot)$ and $F_{Y^d}^{-1}(\cdot) = q(d, \cdot)$ denote the corresponding cumulative distribution and density function and the quantile function. The conditional distribution and conditional quantile are denoted by $F_{Y^d}(\cdot|x)$ and $F_{Y^d}^{-1}(\cdot|x) = q(d, x, \cdot)$.

We are interested in characterizing the dependence structure between Y and D eventually conditioning on a set of covariates W in the presence of essential heterogeneity and in the absence of general equilibrium effects. Knowledge of the joint distribution $(Y^d)_{d \in \mathcal{D}}$ or the conditional joint distribution $(Y^{d|x})_{d \in \mathcal{D}}$ would allow to characterize a distribution for the outcome for any possible level of the treatment. When potential outcomes are comonotonic, they can be described as different functions of the same (single) random variable and QTEs are informative on the impact distribution. The potential outcome could be written as $y^d = q(d, u)$, $u \sim \mathcal{U}(0, 1)$, $q(d, u)$ is increasing in u as is refereed in the literature as the *structural quantile function*. If the potential outcomes are not comonotonic, QTEs are informative on the distance between potential outcomes distributions, which may be interesting per se, but not on the impact distribution. We thus concentrate on strategies that focus on QTEs.² In the binary case, QTEs (see Eq. (2)) are defined as the horizontal distance between the distribution function in the presence and in the absence of the treatment [9, 15].³

$$\delta(\tau) = F_{Y^1}^{-1}(\tau) - F_{Y^0}^{-1}(\tau) \quad 0 < \tau < 1 \quad (2)$$

We can distinguish conditional and unconditional QTEs by characterizing the uniformly distributed random variable that describes the quantile of the outcome variable. This distinction becomes clearer if we think about a specific empirical example.

2.1 Motivating Example: Returns to Education or Training

There is a large literature that studies the returns to education. Key questions in this literature (e.g. does additional education cause wage increase? does additional

²The review will not cover strategies that focus on other objects and may deliver QTEs as by-product such as [8], for instance.

³In the continuous case $\delta(\tau)$ represents the change in Y induced by a change in D from d to $d + \epsilon$ when ϵ is small.

$$\delta(\tau) = \frac{\partial Q_Y(\tau|d)}{\partial d} \quad 0 < \tau < 1 \quad (1)$$

schooling increase wages more for the more able than for the less able? does additional schooling increase or decrease wage inequality?) can be addressed using quantile regression methods. In this applications, the treatment is likely endogenous in the outcome equation without conditioning on additional covariates: typically researchers seek instruments that allow to isolate exogenous variation in education in the wage equation. Suppose we could measure the individual ability a_i that drives the endogeneity of education in the wage equation. Now, consider the alternative specifications for the wage model presented in Eqs. (3) and (4) where D denotes schooling (the treatment).

$$Y_i = \alpha_0(f(\varepsilon_i, a_i)) + \alpha_1(f(\varepsilon_i, a_i))D \quad (3)$$

$$Y_i = \beta_0(\varepsilon_i)a_i + \beta_1(\varepsilon_i)D \quad (4)$$

These specifications differ because they impose different structures of the variables governing the heterogeneity in the returns to education. In Eq. (3) the relative position of an individual in the wage distribution is determined by (ε_i, a_i) , i.e. by both an *unobserved* uniformly distributed error component ε_i and by the *observed* individual ability level while in Eq. (4) the relative position of the individual is only determined by ε_i . In both cases, we can think about the relative position of an individual in the wage distribution as his/her *prone ness* [9] to earn a high wage for a given level of schooling D . However, in model (3) we would refer to the *total prone ness/ability* while in model (4) we would be speaking only about *unobserved prone ness/ability*.⁴ Using model (3) we can explore whether the returns to education vary depending on the individuals' *total* ability levels while using model (4) we can study how the returns to education vary for given observed ability levels. Individuals who earn high wages conditional on some specific level of ability may not be the same individuals who earn high wages in the sample. However conditioning on observed ability maybe important to be able to isolate the causal effect of schooling D on the distribution of wages Y . Equations (5) and (6) represent the structural quantile function corresponding to model (3) and (4), respectively⁵: Eq. (5) is an example of an *unconditional* quantile regression model while Eq. (6) is an example of a *conditional* quantile regression model. This distinction might be empirically relevant since, in general, for a given $\tau \in (0, 1)$, $\alpha_1(\tau) \neq \beta_1(\tau)$.

$$f(\varepsilon, a) \equiv \varepsilon^*, \varepsilon^* \sim \mathcal{U}(0, 1) \quad Q_Y(\tau|d) = \alpha_0(\tau) + \alpha_1(\tau)d \quad (5)$$

$$\varepsilon \sim \mathcal{U}(0, 1) \quad Q_Y(\tau|d) = \beta_0(\tau)a_i + \beta_1(\tau)d \quad (6)$$

⁴To the best of my knowledge, [17] is the first to distinguish between *total* and *observed prone ness*.

⁵Under comonotonicity of potential outcomes, the structural quantile function describes the link between potential outcomes.

3 Identification Strategies and Estimation

In cross-sectional applications, two main identification approaches have been extended to QTEs: strategies based on the *unconfoundedness assumption* and strategies based on the availability of an instrumental variable. In the first case, the researcher must be willing to assume that the joint distribution of the potential outcomes is independent of the treatment conditional on a set of exogenous covariates. Under this assumptions, conditional QTEs can be estimated as originally proposed by Koenker and Bassett [14] and unconditional QTEs can be estimated as proposed by Firpo [10]. Abadie et al. [2] and Chernozhukov and Hansen [6, 7] propose identifying assumptions for conditional quantiles when an instrumental variable is available. The assumptions of Abadie et al. [2] guarantee identification of conditional and unconditional QTEs when the treatment is binary and endogenous and a binary instrument is available. They lead to the moment conditions described in Table 1: in both cases, identification relies on previous results [1, 13] that guarantee that in the sub-population of *compliers* comparisons by treatment D , conditional on X , have a causal interpretation. Recall that *compliers* are individuals whose treatment status is affected by the instrument Z but that this sub-population cannot be identified directly from the data, because it is defined by means of potential outcomes. The moment conditions highlight that is possible to construct weights that “find *compliers* in the population in an average sense” [1]. The weights will differ when one is interested in the conditional or in the unconditional quantiles. Only the weights considered in the second case “simultaneously balance the distribution of the covariates between treated and non-treated *compliers*” [12]. In both cases weights are functions of $P(Z = 1|X)$ and observed variables. Estimation thus proceeds in two steps: (1) weights are estimated; (2) weighted quantile regressions are run.⁶ Estimation requires two steps also under the identification strategy proposed by Chernozhukov and Hansen [6, 7] and [17, 18] but does not involve re-weighting. The crucial assumption for identification in the approach

Table 1 Moment conditions under assumptions in [2] and [11]

Quantile	Conditional	Unconditional
Y_1	$E[\{1(Y < q(1, x)) - \tau\} \cdot w_{y,d,x} \cdot D X] = 0$	$E[\{1(Y < q(1)) - \tau\} \cdot w_{y,d} \cdot D] = 0$
Y_0	$E[\{1(Y < q(0, x)) - \tau\} \cdot w_{y,d,x} \cdot (1 - D) X] = 0$	$E[\{1(Y < q(0)) - \tau\} \cdot w_{y,d} \cdot (1 - D)] = 0$
Weight	$1 - \frac{D[1 - P(Z=1 Y,D,X)]}{1 - P(Z=1 X)} - \frac{(1 - D)P(Z=1 Y,D,X)}{P(Z=1 X)}$	$E[\frac{Z - P(Z=1 X)}{P(Z=1 X)[1 - P(Z=1 X)]} Y, D](2D - 1)$

Note: Positive weights are reported. See [2] and [11] for other definitions of weights

⁶When identification is achieved relying on *unconfoundedness*, the moment conditions are similar but the weights are identically 1 for conditional quantiles [14] and are $\frac{D}{P(D=1|X)} + \frac{1-D}{1-P(D=1|X)}$ for unconditional quantiles [10].

by Chernozhukov and Hansen [6] is *rank invariance* or *rank similarity*, i.e. we require that the individual's rank in the potential outcome distribution, conditional on exogenous covariates, is not systematically affected by the treatment. The assumptions by Chernozhukov and Hansen [6] lead to the moment condition in Eq. (7). Equation (7) suggests an estimation procedure that first requires to compute the conditional quantiles of the random variable $Y - q(d, x, \tau)$ given X and Z ; then, choose as estimate of $q(d, x, \tau)$ the one that minimizes the absolute value of the coefficient associated with Z in the first step.⁷

$$\Pr[Y - q(d, x, \tau) \leq 0 | X, Z] = \tau. \quad (7)$$

The instrumental variable approach for the identification of unconditional QTEs proposed by Powell [17] delivers the moment conditions in Eq. (8)

$$\begin{aligned} E[Z\{1(Y \leq q(d, \tau)) - \tau_X\}] &= 0, \tau_X \equiv P[Y \leq q(d, \tau) | X]. \\ E[1(Y \leq q(d, \tau)) - \tau] &= 0. \end{aligned} \quad (8)$$

These moment conditions reflect the idea that, first, the instrument Z does not affect the distribution of the disturbance once X is controlled for and, second, the joint distribution of X and the disturbance is unrestricted. Estimation involves first an estimation of the quantiles of $Y - q(d, \tau)$ given X and Z and τ_X ; then, a second step choose as estimate of $q(d, \tau)$ the value that minimizes the coefficient of Z averaging over all possible values of X .

We now apply these strategies to two illustrative examples taken from the literature. Table 2 reports estimates of the effect of training (or education) on the conditional and unconditional distribution of earnings (or log wages) using data of males from [2] and data from [5], respectively.⁸ Columns (1) and (2) report results delivered when training or education are treated as exogenous in the estimation of conditional and unconditional quantiles, respectively. Columns (3) and (4) report estimates that address the endogeneity of training or education in the outcome equation relying on [2]. These estimates apply to the sub-population of *compliers*. Columns (5)–(8) report estimates based on [6] or [17]. These approaches guarantee global identification of conditional and unconditional QTEs. We discuss the top-panel estimates first: in the example from [2] the treatment assignment is

⁷This approach can be used when the treatment and instrument are binary, discrete as well as continuous.

⁸In the second example, only reforms that increased compulsory schooling for 3 years are considered (i.e. only Greece, Italy and Finland) and the original treatment (years of education) and instrument (years of compulsory schooling) were recoded to binary. Estimates of columns (1)–(4) have been computed by the author using the STATA package `ivqte` by Froelich and Melly [12], except column (3) for the first example (taken from the article). Estimates in column (1) replicate original results in the papers except that standard errors are now robust to heteroskedasticity; estimates of columns (5)–(8) are taken from [18] for the AAI02 example and obtained using the STATA package `ivqreg` by Do Wan Kwack available from Christian Hansen's research page.

Table 2 Effect of training on the conditional and unconditional distribution of earnings (2], males only) and effect of education on the conditional and unconditional distribution of log wages (5], males, Italy Greece and Finland, treatment and instrument recoded to binary)

Strategy	q	Exogenous training				Endogenous training			
		Conditional		Unconditional		Monotonicity		Rank invariance	
		(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
		KB78	F07	AA102	FM10	CH08	CH08 w/o controls	P11 logit	P11 probit
<i>Effect of training on earnings, Abadie et al. [2] Obs. 5102</i>									
0.25		2510 (417)***	3058 (377)***	702 (670)	414 (754)	530 (629)	200 (746)	100 (753)	100 (750)
0.50		4420 (613)***	4678 (771)***	1544 (1074)	1291 (1239)	310 (1101)	1320 (1234)	790 (1151)	790 (1161)
0.75		4678 (901)***	4626 (1056)***	3131 (1376)**	2457 (1650)	2660 (1845)	1710 (1712)	1490 (1542)	1490 (1530)
0.85		4806 (1045)***	5532 (1241)***	3378 (1811)*	3971 (1886)**	3190 (1185)**	3580 (1427)**	3410 (1542)*	3410 (1550)*
<i>Effect of education on log wages, Brunello et al. [5] Obs. 2292</i>									
0.30		0.168 (0.024)***	0.223 (0.064)***	0.303 (0.142)**	0.514 (16.48)	0.836 (0.063)***	-0.198 (0.033)***	-	-
0.50		0.177 (0.024)***	0.208 (0.062)***	0.328 (0.126)***	0.521 (5.95)	0.985 (0.063)***	-5.119 (0.124)***	-	-
0.75		0.213 (0.026)***	0.297 (0.072)***	0.154 (0.168)***	0.599 (10.13)	1.868 (0.998)**	0.996 (0.037)***	-	-

Column labels refer to the estimation method. KB78: as in [14]; F07: as in [10]; AA102: as in [2]; FM10: as in [11, 12]; CH08: as in [7]; P11: as in [18]
 *p < 0.1, **p < 0.05, ***p < 0.01

randomized thus covariates are not needed for identification. Indeed, under both the identification approaches considered, training effects on the conditional and unconditional quantiles do not exhibit substantial differences in magnitude and all suggest that the effect of training is larger at the top of the earnings distribution.⁹ In addition, both the identification strategies deliver similar results, suggesting that key assumptions are unlikely to be violated in both cases. Let's now turn to the estimates in the bottom part of the table. In this example, covariates are needed for identification: we need to control for country specific secular trends in education and differences across countries in the levels of education and wages to be able to isolate the exogenous variation in education induced by school reforms. In this example, addressing endogeneity seems to have relevant consequences: the estimates in columns (1) and (2) suggest that returns are increasing over the wage distribution, while estimates in column (3) suggest the opposite—although precision of these estimates is low—and in column (4) we find no evidence of heterogeneity.¹⁰ Estimates of conditional QTEs under rank invariance are reported in column (5); estimates of unconditional QTEs in column (6) assume rank invariance and do not use covariates for identification. The estimates in column (5) are unrealistic and suggest that rank invariance is unlikely to hold. Estimates in column (6) are negative and confirm that controlling for covariates is necessary for identification.

4 Conclusions

In this paper, I reviewed approaches that guarantee the identification of QTEs. In many cases, these approaches correspond to extensions of strategies conventionally used in linear regression models (selection on observables, instrumental variables, fixed effects) to quantile regressions. An important consequence of the difference between the statistical tools applied in these two settings is that the interpretation of treatment parameters differs between conditional and unconditional quantile regressions, while, conversely, the law of iterated expectations guarantees that the treatment parameter in a linear regression has both a conditional and an unconditional mean interpretation. It is crucial to bear this in mind while using QTEs to answer a specific research question. Consider the recent proposal of Barlevy and Neal [3] to link educator compensation to the ranks of their students within what the authors call an *appropriately defined comparison sets*. The authors suggest to employ methods in [4] to contrast actual ranks of students of a given teacher with some predicted counterfactual rank. Betebenner [4] however employs *conditional*

⁹When endogeneity of training is addressed, point estimates of the returns to training are generally lower in the unconditional distribution with respect to the returns observed holding race, age, education and marital status fixed.

¹⁰In this example, we look at the effect of three more additional years of schooling on wages. Assuming linearity and dividing point estimates reported by three, the results in columns (1)–(3) are fairly consistent with the literature: association is lower than causal effects; causal estimate suggests a return between 10% and 4% for each additional year of education.

quantile regression methods aimed specifically at answering questions like *Are there students with unusually low growth who need special attention?*, i.e. a value-added specification of the achievement. Barlevy and Neal [3] instead look for a method that allows to isolate the teachers contribution to a student rank in the achievement distribution in a given period, eventually conditioning on covariates for identification. In other words Barlevy and Neal would like to avoid attributing to a teacher the changes in performance of a student that are only due to his initial proficiency level. Standard value-added specifications for students' achievement in quantile regression context are not the appropriate instrument to address questions about the heterogeneity in students' achievement depending on their initial ability level. Those quantile regression describes instead how students experiencing the largest gains in performance over a given time period perform relative to students experiencing the lowest gains in the same period. Cross-sectionally, some of the *high-gain* students may be in the lower part of the test score distribution.¹¹

Acknowledgements This paper benefited from comments by E. Rettore, B. Pacini and F. Mealli. Financial support of MIUR- FIRB 2008 project RBFR089QQC-003-J31J10000060001 grant is gratefully acknowledged.

References

1. Abadie, A.: Semiparametric instrumental variable estimation of treatment response models. *J. Econ.* **113**, 231–263 (2003)
2. Abadie, A., Angrist, J., Imbens, G.: Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings. *Econometrica* **70**, 91–117 (2002)
3. Barlevy, G., Neal, D.: Pay for percentile. NBER Working Paper No. 17194 (2010)
4. Betebenner, D.W.: Norm and criterion-referenced student growth. *Educ. Meas. Issues Pract.* **28**(4), 42–51 (2009)
5. Brunello, G., Fort, M., Weber, G.: Changes in compulsory schooling, education and the distribution of wages in Europe. *Econ. J.* **119**(536), 516–539 (2009)
6. Chernozhukov, V., Hansen, C.: An IV model of quantile treatment effects. *Econometrica* **73**, 245–261 (2005)
7. Chernozhukov, V., Hansen, C.: Instrumental variable quantile regression: a Robust Inference Approach. *J. Econ.* **142**(1), 379–398 (2008)
8. Chesher, A.: Identification in nonseparable models. *Econometrica* **71**, 1405–1441 (2003)
9. Doksum, K.: Empirical probability plots and statistical inference for nonlinear models in the two-sample case. *Ann. Stat.* **2**, 267–277 (1974)
10. Firpo, S.: Efficient semiparametric estimation of quantile treatment effect. *Econometrica* **75**, 259–276 (2007)
11. Froelich, M., Melly, B.: Unconditional quantile treatment effects under endogeneity. IZA Discussion Papers No. 3288 (2010)
12. Froelich, M., Melly, B.: Estimation of quantile treatment effects with STATA. *Stata J.* **10**(3), 423–457 (2010)

¹¹A similar point was made by Powell [16] in his discussion of the analysis of the effect of vouchers on student achievements.

13. Imbens, G., Rubin, D.: Estimating the outcome distribution for compliers in instrumental variables models. *Rev. Econ. Stud.* **64**, 555–574 (1997)
14. Koenker, R., Bassett, G.S.: Regression quantiles. *Econometrica* **46**, 33–50 (1978)
15. Lehmann, E.H.: *Nonparametrics: Statistical Models Based on Ranks*. Holden Day, San Francisco (1974)
16. Powell, D.: Unconditional quantile regression for panel data with exogenous or endogenous treatment variables. RAND Working Paper No. WR-710 (2010)
17. Powell, D.: Unconditional quantile treatment effects in the presence of covariates. RAND Working Paper No. WR-816 (2010)
18. Powell, D.: Unconditional quantile regression for exogenous or endogenous treatment variables. RAND Working Paper No. WR-824 (2011)

Some New Results on the Beta Skew-Normal Distribution

Valentina Marneli and Monica Musio

Abstract

In this paper we study the Beta skew-normal distribution introduced by Marneli and Musio (2013). Some new properties of this distribution are derived including formulae for moments in particular cases and bi-modality properties. Furthermore, we provide expansions for its distribution and density functions. Bounds for the moments and the variance of the Beta skew-normal are derived. Some of the results presented in this work can be extended to the entire family of the Beta-generated distribution introduced by Jones (Test 13(1):1–43, 2004).

1 Introduction

The literature related to the skew-normal distribution (*SN*), introduced in [1], has grown rapidly in recent years. Marneli and Musio [12] proposed a generalization of the skew-normal called Beta skew-normal (*BSN*). This distribution arises quite naturally if we consider the distribution function of the order statistics of the skew-normal distribution and it can also be seen as a special case of the Beta-generated family proposed by Jones [8]. In [12] the authors studied some properties of the Beta skew-normal distribution. In particular, they derived the moment generating function, recurrence relations for moments and two methods for simulating.

The main aim of this paper is to study some new properties of the Beta skew-normal distribution. Particularly, inspired by Gupta and Nadarajah [6], we obtain general expressions for the moments of the *BSN*. Expansions for its distribution and density functions are also provided. Moreover, motivated by the well-known bounds

V. Marneli (✉) • M. Musio
University of Cagliari, via Ospedale 72, Cagliari, Italy
e-mail: marneli.valentina@virgilio.it; mmusio@unica.it

for the moments [5, 7], and the variance of the order statistics [13], we study the problem of finding bounds for the moments and the variance of the Beta-generated family. Some numerical calculations are given in the special case of the BSN. The paper is organized as follows. In Sect. 2 we recall the Beta skew-normal distribution. In Sect. 3 we give expansions for the distribution and the density functions and some theorems about the moments. Furthermore, we provide bounds for the moments and the variance of the Beta-generated distribution. In addition, we give some bimodality properties. All computations have been done using the software R.

2 The Beta Skew-Normal

In this section we remind the BSN class and some of its properties (see [12] for other properties).

Definition 1 A random variable X is said to have a Beta skew-normal distribution, if its density is given by

$$g_{\Phi(x;\lambda)}^B(x; \lambda, a, b) = \frac{2}{B(a, b)} (\Phi(x; \lambda))^{a-1} (1 - \Phi(x; \lambda))^{b-1} \phi(x) \Phi(\lambda x), \quad x \in \mathbb{R}, \quad (1)$$

where $B(a, b)$ denotes the Beta function and $\Phi(x; \lambda)$ is the distribution function of a random variable with skew-normal distribution. We denote the random variable X by $X \sim \text{BSN}(\lambda, a, b)$.

Figure 1 shows the density of the BSN distribution for various values of the parameters. The densities can be unimodal (Fig. 1(a)) or bimodal (Fig. 1(b)).

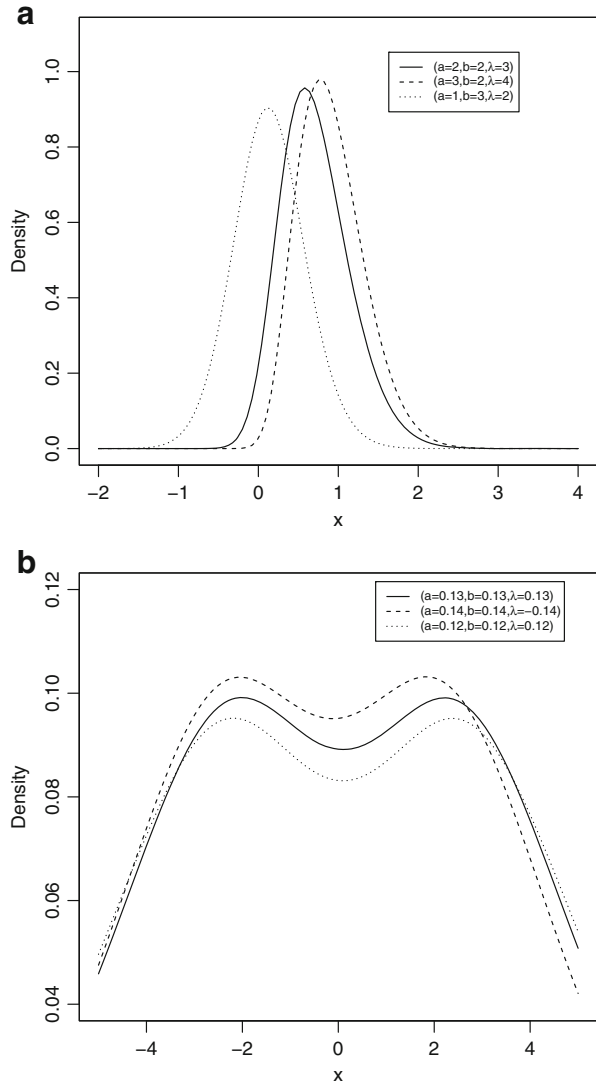
The following properties of the $\text{BSN}(\lambda, a, b)$ have been derived directly from (1).

Properties of $\text{BSN}(\lambda, a, b)$:

- (a) If $a = 1$ and $b = 1$, then we obtain the skew-normal distribution ($\text{SN}(\lambda)$).
- (b) If $\lambda = 0$, then the *BSN* distribution reduces to the Beta-normal one ($\text{BN}(a, b)$) [3].
- (c) If $X \sim \text{BSN}(0, 1, 1)$, then X is a standard normal random variable.
- (d) If $X \sim \text{BSN}(1, \frac{1}{2}, 1)$, then X is a standard normal random variable.
- (e) If $X \sim \text{BSN}(-1, 1, \frac{1}{2})$, then X is a standard normal random variable.
- (f) If $X \sim \text{BSN}(\lambda, a, b)$, then $-X \sim \text{BSN}(-\lambda, b, a)$.
- (g) If $X \sim \text{BSN}(\lambda, a, b)$, then $Y = \Phi(X; \lambda) \sim \text{Beta}(a, b)$ and $Z = 1 - \Phi(X; \lambda) \sim \text{Beta}(b, a)$.
- (h) As $\lambda \rightarrow +\infty$, the *BSN* density tends to the Beta half-normal density [14].

Remark 1 Properties from (c) to (e) establish that the family of $\text{BSN}(\lambda, a, b)$ contains the standard normal distribution as a special case under three different parameter sets. Consequently, we have that a probabilistic model based on the *BSN* distribution is not identifiable under the null hypothesis of normality. It is

Fig. 1 The *BSN* density function for some values of a , b and λ



well known that the classical asymptotic results concerning the likelihood ratio test statistic (LRT) are not true in case of loss of identifiability and the distribution of the LRT statistic is difficult to characterize (see, for example, [10]). A way to avoid this problem is to redefine the model over a modified parameter space such that the null hypothesis of normality is described only by a unique set of parameters.

Moments of the *BSN* have not been evaluated exactly in closed form. However, they have been computed numerically. We have noted that (see [12]):

- for fixed values of a and b the mean and skewness are both increasing function of λ ;
- for fixed values of b and λ the mean and skewness are both increasing function of a ;
- for fixed values of a and λ the mean is a decreasing function of b .

It is also interesting to remind the following results.

Theorem 1 *Let $X \sim \text{BSN}(\lambda, a, b)$ be independent of a random sample (Y_1, \dots, Y_n) from $\text{SN}(\lambda)$, then $X | (Y_{(n)} \leq X) \sim \text{BSN}(\lambda, a + n, b)$ and $X | (Y_{(1)} \geq X) \sim \text{BSN}(\lambda, a, b + n)$, where $Y_{(n)}$ and $Y_{(1)}$ are the largest and the smallest order statistics, respectively.*

Theorem 2 *Let $X \sim \text{BSN}(\lambda, a, b)$ be independent of $Y \sim \text{BSN}(\lambda, c, 1)$ and of $Z \sim \text{BSN}(\lambda, 1, d)$, where c and d are positive real numbers. Then $X | Y \leq X \sim \text{BSN}(\lambda, a + c, b)$ and $X | Z \geq X \sim \text{BSN}(\lambda, a, b + d)$, where $c, d \in \mathbb{R}$.*

Theorem 3 *If $X \sim \text{BSN}(\lambda, a, b)$ is independent of $U_1, \dots, U_n, V_1, \dots, V_m$ then $X | (U_{(n)} \leq X, V_{(1)} \geq X) \sim \text{BSN}(\lambda, a + n, b + m)$, where $U_{(n)} = \max(U_1, \dots, U_n)$ and $V_{(1)} = \min(V_1, \dots, V_m)$.*

For the proof of these theorems we refer to [12]. It is interesting to note that:

Remark 2 Theorem 1 can be used to generate $X \sim \text{BSN}(\lambda, n, 1)$ (see [12]). The $\text{BSN}(\lambda, n, 1)$ can also be generated as follows: let U_1, U_2, \dots, U_n be a random sample from $\text{SN}(\lambda)$, then the random variable $X = \max(U_1, U_2, \dots, U_n)$ has a $\text{BSN}(\lambda, n, 1)$ distribution. How it has been shown in [11], and pointed out in the conclusions of Mameli and Musio [12], the previous theorems can be extended more generally to the class of Jones' Family.

3 New Results on the Beta Skew-Normal Distribution

3.1 Expansion for the Density Function

Here, we give a simple expansion for the BSN density function. By applying the binomial expansion to the distribution function of the $\text{BSN}(\lambda, a, b)$, if b is a real non-integer, we get

$$G_{\Phi(x; \lambda)}^B(x; a, b) = \sum_{i=0}^{\infty} w_i(a, b) \Phi(x; \lambda)^{a+i},$$

where $w_i(a, b) = \frac{1}{B(a, b)} \frac{(b-1)!}{i!(b-1-i)!} \frac{(-1)^i}{a+i}$. Correspondingly, the density function of the BSN can be written as

$$g_{\Phi(x;\lambda)}^B(x; \lambda, a, b) = \sum_{i=0}^{\infty} w_i(a, b) g_{\Phi(x;\lambda)}^B(x; \lambda, a + i, 1), \tag{2}$$

where the weights $w_i(a, b)$ are such that $\sum_{i=0}^{\infty} w_i(a, b) = 1$. However, it is clear from the last equation that $g_{\Phi(x;\lambda)}^B(x; \lambda, a, b)$ can be expressed as an infinite mixture of BSN($\lambda, a + i, 1$) densities with constant weights $w_i(a, b)$. For b integer, the previous sums stop at $b - 1$.

Remark 3 The density function of the SN(λ) distribution can be represented in the following way:

$$\begin{aligned} \phi(z; \lambda) &= 2\phi(z)\Phi(\lambda z) = 2\phi(z)\Phi(\lambda z) (1 - \Phi(z; \lambda) + \Phi(z; \lambda)) \\ &= 2\phi(z)\Phi(\lambda z) (1 - \Phi(z; \lambda)) + 2\phi(z)\Phi(\lambda z)\Phi(z; \lambda) \\ &= \frac{1}{2} \left(g_{\Phi(z;\lambda)}^B(z; \lambda, 1, 2) + g_{\Phi(z;\lambda)}^B(z; \lambda, 2, 1) \right). \end{aligned}$$

In other words the density function of the skew-normal with parameter λ is a mixture between a BSN density with parameters $\lambda, a = 1$ and $b = 2$ and a BSN density with parameters $\lambda, a = 2$ and $b = 1$, which are the density function of the smallest and the largest statistic from a sample of size 2 of a skew-normal distribution with parameter λ , respectively. In general, we can see the density function of the skew-normal with parameter λ as a mixture of Beta skew-normal distributions with the same parameter λ in the following way:

$$\phi(x; \lambda) = \frac{1}{b} g_{\Phi(x;\lambda)}^B(x; \lambda, 1, b) - \sum_{i=1}^{\infty} \frac{(b-1)!}{i!(b-1-i)!} \frac{(-1)^i}{1+i} g_{\Phi(x;\lambda)}^B(x; \lambda, 1+i, 1).$$

The above formula is obtained on setting $a = 1$ in (2) and using the property (a) of the BSN.

We use the previous expansion (2) to present a formula for the moments of the BSN when a and b are integers.

Theorem 4 Let $X \sim \text{BSN}(\mu, \sigma, \lambda, a, b)$ for integer values of a and b , then

$$\begin{aligned} E(X^n) &= \mu^n + \frac{2\mu^n}{B(a, b)} \sum_{j=0}^{b-1} \frac{(-1)^j (b-1)!}{j!(b-1-j)!} \sum_{i=1}^n \frac{(n)!}{i!(n-i)!} \left(\frac{\sigma}{\mu} \right)^i \\ &\quad \times \left\{ \sum_{k=0}^{a+j-1} \frac{(-1)^k (a+j-1)!}{k!(a+j-1-k)!} J_{i,k,\lambda} + (-1)^i J_{i,a+j-1,-\lambda} \right\}, \end{aligned}$$

where

$$J_{i,k,\lambda} = \int_0^\infty z^i \phi(z) \Phi(\lambda z) (1 - \Phi(z; \lambda))^k dz.$$

Proof The proof follows the same lines of that of Theorem 1 in [6].

Remark 4 Clearly, this theorem when $\lambda = 0$ reduces to Theorem 1 in [6]. Furthermore, we can note that the authors, in the cited theorem, defined the function

$$I_{i,k} = \int_0^\infty z^i \phi(z) (1 - \Phi(z))^k dz,$$

which is related to the function $J_{i,k,\lambda}$, when $\lambda = 0$, by the relation $J_{i,k,0} = \frac{1}{2} I_{i,k}$.

The BSN with Parameters 1, n and b (BSN(1, n , b))

As previously noted, general expansions for the moment generating function and the k -th moment of a variable with Beta skew-normal distribution are difficult to find. Exact closed form expressions for the moments can be obtained in certain special cases. One of these cases is discussed in this section.

Theorem 5 *The moment generating function of $X \sim \text{BSN}(1, n, b)$ is*

$$M_X(t) = \frac{2}{B(n, b)} \sum_{j=0}^{\infty} (-1)^j \frac{(b-1)!}{j!(b-1-j)!} e^{\frac{t^2}{2}} E(\Phi^{2(n+j)-1}(V)), \quad (3)$$

where $V \sim N(t, 1)$.

Proof The proof is based on the binomial expansion and the well-known property of the distribution function of the SN distribution $\Phi(x; 1) = \Phi(x)^2$ (see [1]).

We can obtain the moments of $X \sim \text{BSN}(1, n, b)$ readily from the derivatives of $M_X(t)$ in (3). For example, we get the first moment as

$$E(X) = \frac{1}{B(n, b)} \sum_{j=0}^{\infty} \frac{(-1)^j (b-1)!}{j!(b-1-j)!} \frac{1}{n+j} \frac{(2(n+j)-1)(n+j)}{\sqrt{\pi} c_{(2(n+j)-2)} \left(\frac{1}{\sqrt{2}} \right)},$$

where the function $c_m(\beta)$ with $m \in N$ and $\beta \in R$ is the normalizing constant of the Balakrishnan skew-normal density (hereafter denoted by $\text{SNB}_m(\beta)$) (see [16]).

Remark 5 As noted in [16], for $m \geq 4$, there is no closed form for $c_m(\beta)$, but one can find some approximate values for $c_m(\beta)$ using Table 1 in [17]. Furthermore, we can note that in the special case $b = 1$ and $n = 2$, we have $E(X) = \frac{6}{\sqrt{\pi}} \left[\arctan(\sqrt{2}) \right]$, which is exactly the mean of the maximum from a sample of size 2 from a SN(1) obtained in [2].

The following theorem provides a recursive formula for the moments of the BSN(1, n , b):

Theorem 6 *Let $X \sim \text{BSN}(1, n, b)$. Then*

$$E(X^k) = \sum_{j=0}^{\infty} \frac{(-1)^j (b-1)!}{B(n, b) j! (b-1-j)! (n+j)} \times \left\{ (k-1)E(Y^{k-2}) + \frac{2n+2j-1}{2^{\frac{k+1}{2}} \sqrt{\pi}} c_{(2(n+j)-2)}\left(\frac{1}{\sqrt{2}}\right) E(W^{k-1}) \right\},$$

where $W \sim \text{SNB}_{(2(n+j)-2)}\left(\frac{1}{\sqrt{2}}\right)$ and $Y \sim \text{SNB}_{(2(n+j)-1)}(1)$.

Proof The proof follows by combining the recursive formula of the moments of the Balakrishnan skew-normal distribution [16] with the expansion

$$E(X^k) = \frac{2}{B(n, b)} \sum_{j=0}^{\infty} \frac{(-1)^j (b-1)!}{j! (b-1-j)! 2(n+j)} E(Y^k),$$

where $Y \sim \text{SNB}_{(2(n+j)-1)}(1)$.

Remark 6 It should be noted that similar results can be also proved for the BSN(-1, b , n) distribution. This is due to the fact that, as stated in property (f) of Sect. 2, if $X \sim \text{BSN}(1, n, b)$ then $-X \sim \text{BSN}(-1, b, n)$.

3.2 Bounds of the Moments and the Variance of the Beta – F

Several authors have given methods of finding bounds for the moments of order statistics. One of the earliest result is that derived in [5, 7]. Different methods are required for the variance of the order statistics. Following the idea of these works, we apply Hölder's inequality and Hoeddfing's identity to find inequalities for the moments and the variance of the Beta-generated distribution.

Bounds of the Moments

In this section we assume that X and Y have distributions $G_{F(\cdot)}^B$ and $F(\cdot)$, respectively.

Theorem 7 Let $k > 0$, $p > 1$ and $E(Y^{kp}) < \infty$. Then we have

$$E(X^k) \leq \frac{1}{B(a, b)} (E(Y^{kp}))^{\frac{1}{p}} \left(B\left(\frac{pa-1}{p-1}, \frac{pb-1}{p-1}\right) \right)^{1-\frac{1}{p}}.$$

Proof Proof is based on Hölder's inequality.

It is worth pointing out that the upper bound of the k -th moment of a random variable with a *Beta* – F distribution depends on the kp -th moment of a r.v. with a F distribution. In the specific case of the BSN, the k -th moment of a r.v. with a Beta skew-normal distribution can be bounded above by a function of the kp -th moment of a r.v. with a skew-normal distribution.

Bounds of the Variance of the Beta-Generated Distribution

Let $X \sim G_{F(\cdot)}^B(\cdot, a, b)$, with $a > 1$ and $b > 1$. We are interested in finding a bound for the variance of X as a function of the variance of $Y \sim F(\cdot)$. Let us introduce the notations:

$$\begin{aligned} G(x) &= I_x(a, b), & g(x) &= G'(x), & t_1(x) &= \frac{G(x)}{x}, \\ t_2(y) &= \frac{1 - G(y)}{1 - y}, & t(x, y) &= t_1(x)t_2(y), & t(x) &= t(x, x), \end{aligned}$$

with $0 < x \leq y < 1$ and $I_x(a, b)$ denotes the incomplete Beta ratio. The following lemma is a generalization of the Lemma 2.1 in [13].

Lemma 1 Let $a > 1$ and $b > 1$. Then there exist unique numbers $\rho_1 = \rho_1(a, b)$, $\rho_2 = \rho_2(a, b)$ satisfying $0 < \rho_1 < \frac{a-1}{a+b-2} < \rho_2 < 1$, such that, for $0 < x < y < 1$:

1. $t_1(x)$ strictly increases in $(0, \rho_2)$ and strictly decreases in $(\rho_2, 1)$ and similarly $t_2(y)$ strictly increases in $(0, \rho_1)$ and strictly decreases in $(\rho_1, 1)$.
2. If $x \geq \rho_1$ or $y \leq \rho_2$, then $t(x, y) < \max\{t(x), t(y)\}$.
3. If $x < \rho_1$ and $y > \rho_2$, then $t(x, y) < t(\rho_1, \rho_2) < \max\{\rho_1, \rho_2\}$.
4. There exists a unique $x_0 = x_0(a, b) \in (\rho_1, \rho_2)$ such that the function $t(x)$ strictly increases in $(0, x_0)$ and strictly decreases in $(x_0, 1)$.

Proof The proof follows similar approach to that used in lemma 2.1 in [13].

Definition 2 The variance function $\sigma_b^2(a)$ is defined by the following relation

$$\sigma_b^2(a) = \sup_{0 < x < 1} \left(\frac{G(x)(1 - G(x))}{x(1 - x)} \right), \quad a \geq 1 \text{ and } b \geq 1.$$

Remark 7 It is of interest to point out that $\sigma_b^2(a)$, as $\sigma_n^2(k)$ defined in [13], does not have a closed form. However, it is possible to identify the following behaviour of the function $\sigma_b^2(a)$:

- if $b = 1$, then $\sigma_b^2(a) = a$;
- if $a = b = 1$, then $\sigma_b^2(a) = 1$.

Theorem 8 Let $X \sim G_{F(\cdot)}^B(\cdot; a, b)$, $Y \sim F(\cdot)$, $a \geq 1$ and $b \geq 1$. Then

$$\text{Var}(X) \leq \sigma_b^2(a)\text{Var}(Y). \tag{4}$$

Proof The proof proceeds along lines similar to that of Theorem 3.1 at page 189 in [13], which is based on Hoeffding’s identity for the covariance.

In order to compare the left- and right-hand sides of the inequality (4) when $X \sim \text{BSN}(\lambda, a, b)$, we have carried out a numerical study. Some of the results found are reported in Table 1. The study focuses on some particular values of the parameters a and b , attention is given to the case $a = b = 1$, i.e. when the BSN distribution reduces to a SN distribution and $\sigma_{\text{BSN}}^2 = \text{Var}(X)$ and $\sigma_b^2(a)\text{Var}(Y)$ coincide. Moreover, in such case if $\lambda = 0$, then $\text{Var}(Y) = 1$ and $\sigma_{\text{BSN}}^2 = \text{Var}(X)$ coincides with $\sigma_b^2(a)$. Furthermore, we observe that if the parameter a (resp. b) takes the value 1 and if b (resp. a) is “large”, then $\sigma_b^2(a)\text{Var}(Y)$ is not close to σ_{BSN}^2 .

Table 1 The variance of the BSN(λ, a, b) and $\sigma_b^2(a)\text{Var}(Y)$ for different values of a, b and λ

a	b	λ	σ_{BSN}^2	$\sigma_b^2(a)\text{Var}(Y)$
1	1	-10	0.3696834	0.3696834
1	1	0	1	1
1	1	10	0.3696834	0.3696834
1	10	-10	0.2625293	3.696335
1	10	0	0.3443438	9.99865
1	10	10	0.01859684	3.696335
2	10	-10	0.139195	0.8695345
2	10	0	0.2051976	2.352106
2	10	10	0.02108108	0.8695345
10	1	-10	0.01859684	3.696335
10	1	0	0.3443438	9.99865
10	1	10	0.2625293	3.696335
10	10	-10	0.03145790	0.3696834
10	10	0	0.08079098	1
10	10	10	0.03145790	0.3696834

3.3 Bimodal Properties

The densities of the *BSN* family can be unimodal or bimodal depending on the choice of the parameters [12]. We prove, in this section, some bi-modality properties of the *BSN* distribution.

Theorem 9 *A mode of $BSN(\lambda, a, b)$ is a point $x_0 = x_0(\lambda, a, b)$ that satisfies*

$$x_0 = \left\{ \frac{\lambda \phi(\lambda x_0)}{\Phi(\lambda x_0)} + (a-1) \frac{\phi(x_0; \lambda)}{\Phi(x_0; \lambda)} - (b-1) \frac{\phi(x_0; \lambda)}{1 - \Phi(x_0; \lambda)} \right\} \quad (5)$$

and $\frac{\partial r(x_0)}{\partial x} \leq 0$ where

$$\begin{aligned} r(x) = & -x\Phi(\lambda x)\Phi(x; \lambda)(1 - \Phi(x; \lambda)) + \lambda\phi(\lambda x)\Phi(x; \lambda)(1 - \Phi(x; \lambda)) \\ & + (a-1)\Phi(\lambda x)\phi(x; \lambda)(1 - \Phi(x; \lambda)) - (b-1)\Phi(\lambda x)\phi(x; \lambda)\Phi(x; \lambda). \end{aligned}$$

Proof The proof consists in finding the maximum points of the density (1).

From this theorem we easily derive the following corollaries.

Corollary 1 *If $BSN(\lambda, a, b)$ has a mode at x_0 , then $BSN(-\lambda, b, a)$ has a mode at the point $-x_0$.*

Proof The proof follows similar lines of the one in Proposition 2.3 in [15]. Let $q(x)$ be the function obtained replacing a with b (resp. b with a) and λ with $-\lambda$. By noting that $q(x) = -r(-x)$, if x_0 is a modal point for $BSN(\lambda, a, b)$ then $-x_0$ is a modal point for $BSN(-\lambda, b, a)$.

Corollary 2 *The bimodal property of $BSN(\mu, \sigma, \lambda, a, b)$ is independent of the parameters μ and σ .*

Remark 8 Note that also the modal point of the $BN(\mu, \sigma, a, b)$ is independent of the parameters μ and σ (see the pioneering approach in [4] and the more detailed results in [15]).

3.4 The Square of a *BSN*

Finally, we conclude the section presenting a theorem about the limiting behaviour of the square of a $BSN(\lambda, a, b)$.

Theorem 10 *If $X \sim \text{BSN}(\lambda, a, b)$, then $X^2 \rightarrow B\chi^2(1, a, b)$, as $\lambda \rightarrow \infty$, where $B\chi^2(1, a, b)$ is a Beta chi-square random variable with parameters 1, a and b .*

Proof Let $Y = X^2$. The density of the random variable Y is

$$f_Y(y) = \frac{f_{\chi^2(1)}(y)}{B(a, b)} \left\{ \Phi(\lambda\sqrt{y}) (\Phi(\sqrt{y}; \lambda))^{a-1} (1 - \Phi(\sqrt{y}; \lambda))^{b-1} + \Phi(-\lambda\sqrt{y}) (\Phi(-\sqrt{y}; \lambda))^{a-1} (1 - \Phi(-\sqrt{y}; \lambda))^{b-1} \right\}, \quad y > 0,$$

where $f_{\chi^2(1)}(\cdot)$ is the chi-square density function. As $\lambda \rightarrow \infty$, the term inside curly brackets converges to $(2\Phi(\sqrt{y}) - 1)^{a-1} (2(1 - \Phi(\sqrt{y})))^{b-1}$, where $2\Phi(\sqrt{y}) - 1$ is the distribution function of a chi-square random variable with 1 degree of freedom.

The Beta chi-square distribution is a special case of the Beta–Gamma distribution introduced in [9].

Acknowledgements The authors acknowledge helpful suggestions from the editors and an anonymous referee.

References

1. Azzalini, A.: A class of distributions which includes the normal ones. *Scand. J. Stat.* **12**(2), 171–178 (1985)
2. Chiogna, M.: Some results on the scalar skew-normal distribution. *J. Ital. Stat. Soc.* **7**, 1–13 (1998)
3. Eugene, N., Lee, C., Famoye, F.: Beta-normal distribution and its applications. *Commun. Stat. Theory Methods* **31**(4), 497–512 (2002)
4. Famoye, F., Lee, C., Eugene, N.: Beta-normal distribution: bimodality properties and application. *J. Mod. Appl. Stat. Methods* **3**(1), 85–103 (2004)
5. Gumbel, E.: The maxima of the mean largest value and of the range. *Ann. Math. Stat.* **25**, 76–84 (1954)
6. Gupta, A.K., Nadarajah, S.: On the moments of the beta-normal distribution. *Commun. Stat. Theory Methods* **33**(1), 1–13 (2005)
7. Hartley, H., David, H.: Universal bounds for the mean range and extreme observation. *Ann. Math. Stat.* **25**, 85–99 (1954)
8. Jones, M.C.: Families of distributions arising from distributions of order statistics. *Test* **13**(1), 1–43 (2004)
9. Kong, L., Lee, C., Sepanski, J.H.: On the properties of beta-gamma distribution. *J. Mod. Appl. Stat. Methods* **6**(1), 187–211 (2007)
10. Lindsay, B.G.: *Mixture Models: Theory, Geometry and Applications*. IMS, Hayward (1995)
11. Mameli, V.: Two generalizations of the skew-normal distribution and two variants of McCarthy's Theorem. Doctoral dissertation, Cagliari University, Italy (2012)
12. Mameli, V., Musio, M.: A generalization of the skew-normal distribution: the beta skew-normal. *Commun. Stat. Theory Methods*, **42**(12), 2229–2244 (2013)
13. Papadatos, N. : Maximum variance of order statistics. *Ann. Inst. Stat. Math.* **47**(1), 185–193 (1995)

14. Pescim, R.R., Demétrio, C.G.B., Cordeiro, G.M., Ortega, E.M.M., Urbano, M.R.: The beta generalized half-normal distribution. *Comput. Stat. Data Anal.* **54**(4), 945–957 (2010)
15. Rêgo, L.C., Cintra, R.J., Cordeiro, G.M.: On some properties of the beta normal distribution. *Commun. Stat. Theory Methods* **41**(20), 3722–3738 (2012)
16. Sharafi, M., Behboodan, J.: The Balakrishnan skew-normal density. *Stat. Pap.* **49**, 769–778 (2008)
17. Steck, G.P.: Orthant probabilities for the equicorrelated multivariate normal distribution. *Biometrika* **49**(3–4), 433–445 (1962)

The Median of a Set of Histogram Data

Lidia Rivoli, Rosanna Verde, and Antonio Irpino

Abstract

According to Symbolic Data Analysis, a histogram variable describes each object by means of a histogram of values to an object rather than a single value. In the literature, the definition of the average and the standard deviation has been extended to histogram variables. In this paper, we propose a definition and an algorithm for extracting the main order statistics, the median and quartiles, for a histogram variable observed on a set of units. In particular, for we propose to define a median histogram according to a criterion that minimizes the sum of ℓ_1 Wasserstein distances, a particular probabilistic metric, between distributions. We show that the solution of the problem requires to search for a level-wise order defined on the quantile functions (the inverse of the cumulative distribution functions) of the corresponding histogram data. Evidences from an application on real data show that the proposed order statistics for a histogram variable have similar properties to the classic order statistics for a single-valued variable. Finally, we propose two skewness indices for a histogram variable based on the comparison between the average and the median quantile functions.

L. Rivoli

Department of Mathematics and Statistics, University of Naples Federico II, Naples, Italy
e-mail: lidia.rivoli@unina.it

R. Verde • A. Irpino (✉)

Department of Political Sciences “J.Monnet”, Second University of Naples, Caserta, Italy
e-mail: rosanna.verde@unina2.it; antonio.irpino@unina2.it

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,
Studies in Theoretical and Applied Statistics,
DOI 10.1007/978-3-319-27274-0_4

1 Introduction

Currently, large amount of data are produced from several applicative fields: telecommunications, sensor networks, social networks, transactional systems, to name a few. Often, the objective of an analysis of such data consists in identifying regularities, distributions or, in general, summaries, and then to extract further knowledge from more structured representations with the minimum loss of interesting information. The representation of a group of observations through a histogram provides a synthetic representation about location, scale, and shape of the data distribution.

Symbolic Data Analysis [2,3], a recent approach for the analysis of complex data, provides new tools for the treatment of data described by *Histogram Variables*. A *Histogram Variable* is particular type of *symbolic quantitative multi-valued modal variables* providing a description for each unit by a histogram of values rather than a single numerical one.

In this paper, we propose a definition of histogram-order statistics for a *Histogram Variable* in the framework of the Symbolic Data Analysis (in short SDA). Order statistics are generally related to the definition of an order relation on the descriptions. For example, it is known that the median is the central value of a ordered set of numbers. Histogram data, being set-valued descriptions, cannot be generally ordered. However, we consider that the median is also the value that minimize the sum of ℓ_1 distances from the observed values. We adopt the same principle for proposing the median histogram, and, then for the quartile histograms. We propose to use a suitable ℓ_1 metric chosen among the class of Wasserstein–Kantorovich–Monge–Gini metrics [7]. For example, the ℓ_2 Wasserstein distance has been employed for the definition of the sample mean of a histogram variable [6], and it has shown several interpretative advantages for histogram data [8]. These distances are based on quantile functions (the inverse of the cumulated distribution functions) associated with their corresponding density functions that, for histograms, are easy to compute.

Following a similar principle, the median histogram is found as the histogram which minimizes the sum of the ℓ_1 Wasserstein distance with respect to all the observed histograms. We show that the proposed method provides a median histogram such that its quantiles are the median of the corresponding quantiles of the observed histograms (thus it is a quantile-wise, or level-wise, order). Further, we extend the same procedure in order to provide the definitions of the quartile histograms of a histogram variable. Finally, we propose a skewness function based on the quantile-wise comparison between the histograms related to the median and the mean quantile functions of a histogram variable.

The paper is organized as follows. In Sect. 2, we introduce the ℓ^p Wasserstein metric used to compare histogram data; in Sect. 3, we furnish a detailed description of algorithm for computing the median and the quartile quantile functions. Section 4 shows the proposal for a skewness function for a histogram variable. Section 5

provides some evidences of the proposed methodologies on real data, while, Sect. 6 ends the paper.

2 Wasserstein Metric for Histogram Data

According to [2] and [5], given a set E of N units and a real variable Y a histogram variable is a mapping $\mathbf{H}_Y : E \rightarrow \mathcal{B}_Y$ where \mathcal{B}_Y is the set of all the possible histograms with domain in Y . The histogram description of the generic unit i ($i = 1, \dots, n$) is a realization H_{iY} of the variable \mathbf{H}_Y and it is defined by a finite number of pairs $\{(I_{ik}, f_{ik}); k = 1, \dots, K_i\}$ where $I_{ik} = [\underline{y}_{ik}, \bar{y}_{ik})$ (with $\underline{y}_{ik} < \bar{y}_{ik}$), is the k -th generic bin of the histogram and $f_{ik} \geq 0$ is the associated relative frequency such that $\sum_{k=1}^{K_i} f_{ik} = 1$. For each bin, we observe the cumulated relative frequencies w_{ik} such that:

$$w_{ik} = \sum_{l=1, \dots, k} f_{il}, \quad k = 1, \dots, K_i \quad (1)$$

According to the definition of histogram, the values in each bin, denoted as $I_{ik} = [\underline{y}_{ik}, \bar{y}_{ik})$, $\forall i = 1, \dots, N$, are uniformly distributed. Therefore, the cumulative distribution function (*cdf*) $F_i(y)$ associated with each H_{iY} is the piecewise linear function defined as follows:

$$F_i(y) = \begin{cases} 0 & \text{if } y < \underline{y}_{i1} \\ w_{ik-1} + \frac{y - \underline{y}_{ik}}{\bar{y}_{ik} - \underline{y}_{ik}} f_{ik} & \text{if } \underline{y}_{ik} \leq y < \bar{y}_{ik}, k = 1, \dots, K_i \\ 1 & \text{if } y \geq \bar{y}_{iK_i} \end{cases} \quad (2)$$

Under few strictly conditions,¹ its inverse, the quantile function *qf*, is the piecewise linear function defined as follows:

$$F_i^{-1}(t) = \begin{cases} \underline{y}_{i1} & \text{if } t = 0 \\ \underline{y}_{ik} + \frac{t - w_{ik-1}}{w_{ik} - w_{ik-1}} (\bar{y}_{ik} - \underline{y}_{ik}) & \text{if } w_{ik-1} \leq t < w_{ik}, k = 1, \dots, K_i \\ \bar{y}_{iK_i} & \text{if } t = 1 \end{cases} \quad (3)$$

Figure 1 shows an example of the *cdf* and *qf* corresponding to a histogram datum. In order to compare two histogram data, a suitable metric is needed. A natural way to choose a metric can be referred to those proposed for comparing probability distributions. In [8] different metrics for histogram data are discussed: the

¹It is worth noting that $F_i(y)$ is invertible between in $[\underline{y}_{i1}; \bar{y}_{iK_i}]$ if and only if $f_{ik} > 0$ for each $k = 1, \dots, K_i$.

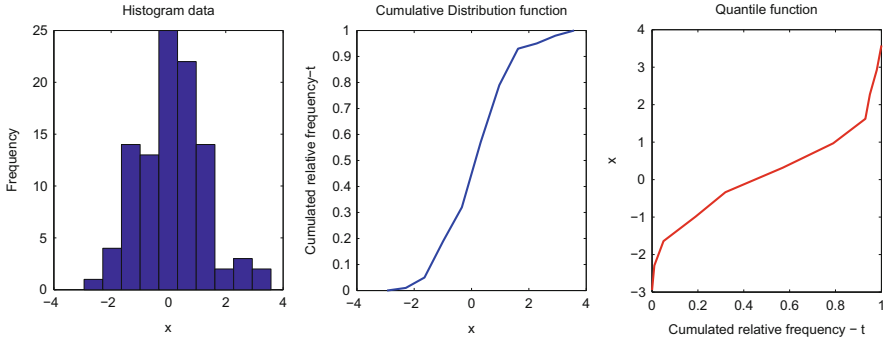


Fig. 1 From the *left to the right*: a histogram datum, its cumulative distribution function (*cdf*) and the corresponding quantile function (*qf*)

f -divergence based measures, the discrepancy metric, the Kolmogorov (or Uniform metric), the Prokhorov–Levi distance, and the ℓ_2 Wasserstein distance.

In this paper, we focus on the ℓ^p Wasserstein distance [7], a distance based on quantile functions, that appears one of the most suitable tool for comparing histogram data, both for its computational and interpretative properties (for a detailed examination of the advantages of the Wasserstein distance see [8] and [9]). In the following, instead of H_{iY} , we denote with H_i the generic i -th histogram datum since we consider a single histogram variable and for providing a clearer notation.

Given two histograms $H_i = \{(I_{ik}, f_{ik}); k = 1 \dots, K_i\}$ and $H_j = \{(I_{jk}, f_{jk}); k = 1 \dots, K_j\}$ and their corresponding F_i^{-1} and F_j^{-1} qf s, the generic ℓ^p Wasserstein distance [7] is defined as follows:

$$d_p(H_i, H_j) = \left\{ \int_0^1 |F_i^{-1}(t) - F_j^{-1}(t)|^p dt \right\}^{\frac{1}{p}}. \quad (4)$$

It is worth noting that a closed form of the distance depends from the possibility of expressing the quantile functions in closed form too. In [9], it is provided a closed form of the (4) for $p = 2$:

$$d_2(H_i, H_j) = \left\{ \int_0^1 (F_i^{-1}(t) - F_j^{-1}(t))^2 dt \right\}^{\frac{1}{2}}. \quad (5)$$

Considering the cumulated frequencies w_{ik} , $k = 1, \dots, K_i$ and w_{jk} , $k = 1, \dots, K_j$ associated to the elementary bins I_{ik} and I_{jk} , the union of two sequences can be expressed as:

$$\{w_{i0}, \dots, w_{iu}, \dots, w_{iK_i}\} \cup \{w_{j0}, \dots, w_{ju}, \dots, w_{jK_j}\}. \quad (6)$$

The sorted values, without repetitions, are represented by the set:

$$w = \{w_0, \dots, w_l, \dots, w_L\},$$

where: $w_0 = 0$, $w_L = 1$, and $\max(K_i, K_j) \leq L \leq (K_i + K_j - 1)$. Each interval (w_{l-1}, w_l) allows us to identify two new uniformly dense bins, one for H_i and one for H_j , having respectively the following bounds:

$$I_{il}^* = [F_i^{-1}(w_{l-1}); F_i^{-1}(w_l)] \quad \text{and} \quad I_{jl}^* = [F_j^{-1}(w_{l-1}); F_j^{-1}(w_l)],$$

with both associated weights equal to $f_l^* = w_l - w_{l-1}$. Each bin I_{il}^* can be expressed as $I_{il}^* = c_{il} + r_{il}(2t - 1)$ for $0 \leq t \leq 1$ where $c_{il} = (F_i^{-1}(w_l) + F_i^{-1}(w_{l-1}))/2$ is the corresponding center (or midpoint) and $r_{il} = (F_i^{-1}(w_l) - F_i^{-1}(w_{l-1}))/2$ is the corresponding radius (or half-length). Thus, the squared distance in (5) can be written as the following sum with a finite number of terms:

$$d_2^2(H_i, H_j) := \sum_{l=1}^m f_l^* \left[(c_{il} - c_{jl})^2 + \frac{1}{3} (r_{il} - r_{jl})^2 \right]. \quad (7)$$

Furthermore, according to [6], it is possible to define the average histogram of a set of N histograms. The average histogram denoted with H_{AV} is the histogram that minimizes the following sum of ℓ_2 Wasserstein distances:

$$\sum_{i=1}^N d^2(H_i, H_{AV}) = \sum_{i=1}^N \sum_{l=1}^L f_l^* \left[(c_{il} - c_l)^2 + \frac{1}{3} (r_{il} - r_l)^2 \right]. \quad (8)$$

It is easy to prove that the (8) reaches a minimum when:

$$c_l = N^{-1} \sum_{i=1}^N c_{il} \quad ; \quad r_l = N^{-1} \sum_{i=1}^N r_{il}.$$

Finally, the quantile function associated with the average histogram is the average quantile function denoted with $AV(t)$, $0 \leq t \leq 1$.

3 The Median Quantile Function and Median Histogram

In this section, we define the *median quantile function* and the respective *median histogram* for a set of histogram data H_i ($i = 1, \dots, N$).

According to [1], the *median histogram* is the histogram H_{ME} that minimizes the sum of the ℓ_1 Wasserstein distances from all the observed histograms. It is obtained

from the following minimization problem:

$$H_{\text{ME}} = \arg \min_H \sum_{i=1}^N d_1(H_i, H) \Leftrightarrow F_{\text{ME}}^{-1}(t) = \arg \min_{F^{-1}(t)} \sum_{i=1}^N \int_0^1 |F_i^{-1}(t) - F^{-1}(t)| dt, \quad (9)$$

where F_i^{-1} and F_{ME}^{-1} are the quantile functions associated to H_i and H_{ME} , respectively.

Similarly to the standard statistical variables, the solution is unique when N is odd, while in the even case, the definition of a median histogram is not unique [1]. Analogously to the classic case, we overcome this limitation using the average (computed with respect to the ℓ_2 Wasserstein distance) between the two most central quantile functions.

Being the distance designed for continuous functions, we show that the solution of the minimization problem (9) can be obtained in a finite time and that produces a median quantile function that is level-wise central. In the following, we present the algorithm for obtaining the *level-wise Median quantile function*, namely, a quantile function that, at each $t \in [0, 1]$, leaves, in the odd case, $(N - 1)/2$ quantiles before and $(N - 1)/2$ after it. We remark that we reach a solution that is consistent with a level-wise order among the quantile functions (the quantile functions are ordered for each level of t) but it is not generally consistent with a full order or semi-order relation among quantile functions. Naturally, if for each $t \in [0, 1]$ the order of the quantile functions is always the same, the level-wise order can be interpreted as full order relation but, in practice, this is infrequent (quantile functions usually intersect each other).

The algorithm is performed in two consecutive steps: a homogenization step and a selection step.

3.1 Homogenization Step

The homogenization step is required for finding a partition of the domain of the quantile functions $[0; 1]$ into intervals of levels of cumulated relative frequencies such that for each interval of levels all the quantile functions are smooth. For this aim, we consider the set containing all the cumulated relative frequencies associated with the N histograms H_i ($i = 1, \dots, N$) as follows:

$$\{w_{10}, \dots, w_{1K_1}, \dots, w_{i1}, \dots, w_{iK_i}, \dots, w_{N0}, \dots, w_{NK_N}\}, \quad (10)$$

and then we obtain the set w of sorted and unique cumulated relative frequencies denoted as follows:

$$w = \{w_0, \dots, w_l, \dots, w_L\}, \quad (11)$$

where $w_0 = 0$, $w_L = 1$. Denoting with $\bar{K} = N^{-1} \sum_{i=1}^N K_i$ the average number of bins, it is straightforward to show that the value of L is bounded as follows:

$$\max_{1 \leq i \leq N} K_i \leq L \leq (N(\bar{K} - 1) + 1).$$

Since the quantile functions are piecewise linear, the values of $F_i^{-1}(w_l)$ for each w_l , $l = 1, \dots, L$ and for each $i = 1, \dots, N$, are computed by a simple linear interpolation. Therefore, each H_i ($i = 1, \dots, N$) is expressed by a new set of L pairs $\{(I_{il}^*, f_{il}^*); l = 1, \dots, L\}$, where: $I_{il}^* = [F_i^{-1}(w_{l-1}), F_i^{-1}(w_l)]$ and $f_{il}^* = w_l - w_{l-1}$.

3.2 Median Level Piecewise Selection Step

This step is repeated for each elementary intervals $[w_{l-1}; w_l]$. After the homogenization step and for all the N histogram data, each interval of levels is the support of the segment starting from the point of coordinates $(w_{l-1}; F_i^{-1}(w_{l-1}))$ and ending at the point $(w_l; F_i^{-1}(w_l))$ (with $w_{l-1} \leq t < w_l$) of the generic i -th quantile functions.

Let $F_{(i)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$) the i -th piece-quantile function, with (i) its order with respect to all the other pieces quantile functions $F_{(j)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$). The order (i) of $F_{(i)}^{-1}(t)$ is kept in the interval $[w_{l-1}, w_l]$ only if there are not intersections with other pieces of quantile functions $F_{(j)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$ and $j = 1, \dots, N, j \neq i$). Indeed, in case of intersections in $[w_{l-1}, w_l]$ the order of the quantile functions may change in some points of the interval (an example is shown in Fig. 2). If intersection points are detected, then a further split of the interval in sub-intervals of levels is performed, the set of levels w_l increases and the final \mathbf{w}

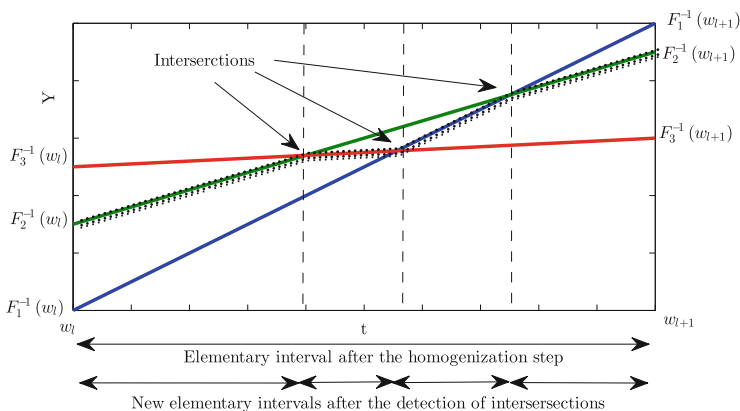


Fig. 2 Partition of an elementary interval $[w_l; w_{(l+1)}]$ (after the homogenization step) into sub-intervals and selection of the new pieces of the median-qf (the dotted path) with $N = 3$

set is updated by the new intervals of levels. Considering the search of the median quantile function, the selection of the piece $F_{\left(\frac{N+1}{2}\right)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$) related to the median quantile function is performed on a higher number of intervals of levels $[w_{l-1}, w_l)$ with $l = 1, \dots, m$ (with $m \geq L$) than the initial one.

When all the pieces belong to the same quantile function, the obtained median quantile function corresponds to an observed quantile function $F_{\left(\frac{N+1}{2}\right)}^{-1}(t)$. However, this rarely occurs in practice. In this case, the median quantile function is obtained by joining the selected $F_{\left(\frac{N+1}{2}\right)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$) pieces that belong to different observed quantile functions for each interval of levels.

We denote the median quantile function as $ME(t)$, in order to distinguish it with respect to an observed quantile function. Finally, the median histogram is the histogram associated with the median quantile function.

Similarly to the classic case, if N is even, $ME(t)$ corresponds to the average quantile function of the two most central quantile functions $F_{\left\lfloor \frac{N}{2} \right\rfloor}^{-1}(t)$ and $F_{\left\lceil \frac{N}{2} \right\rceil}^{-1}(t)$ ²:

$$ME(t) = \frac{F_{\left\lfloor \frac{N}{2} \right\rfloor}^{-1}(t) + F_{\left\lceil \frac{N}{2} \right\rceil}^{-1}(t)}{2} \quad t \in [0; 1] \quad (12)$$

With slight modifications of the algorithm, we may obtain the generic $p \cdot N$ ($p \in [0; 1]$) order-statistic quantile functions denoted with $Q_{(p \cdot N)}(t)$. Similarly, the order-statistic histogram corresponds to the order-statistic quantile functions uniquely: for example, the first-quartile histogram H_{Q_1} is related to the quantile function $Q_{\left(\frac{1}{4} \cdot N\right)}(t)$, the third-quartile histogram H_{Q_3} to the quantile function $Q_{\left(\frac{3}{4} \cdot N\right)}(t)$ and, as said before, the median histogram H_M to the quantile function $ME(t)$.

4 A Skewness Index for Histogram Data

In exploratory data analysis, a skewness index of an empirical distribution can be obtained comparing the average and the median values. In this paper, we propose a skewness measure for a histogram variable based on a normalized

² $\left\lfloor \frac{N}{2} \right\rfloor$ denotes the floor function, while $\left\lceil \frac{N}{2} \right\rceil$ denotes the ceil function.

difference between the average and the median quantile functions. Firstly, a *level-wise skewness index* is proposed as follows:

$$\mathcal{A}_p(l) = \text{sign}(\delta_l) \left\{ \frac{\int_{w_{l-1}}^{w_l} |\text{AV}(t) - \text{ME}(t)|^p dt}{N^{-1} \sum_{i=1}^N \int_0^1 |F_i^{-1}(t) - \text{AV}(t)|^p dt} \right\}^{\frac{1}{p}} \quad (13)$$

where the sign (δ_l) is the sign of the following quantities computed for each interval of levels $[w_{l-1}, w_l]$:

$$\delta_l = \int_{w_{l-1}}^{w_l} (\text{AV}(t) - \text{ME}(t)) dt = \int_{w_{l-1}}^{w_l} \text{AV}(t) dt - \int_{w_{l-1}}^{w_l} \text{ME}(t) dt.$$

We consider the index expressed in (13) only for $p = 1, 2$. In these cases, the normalization factor can be considered as the *mean absolute deviation* and the *standard deviation* [9] of a histogram variable, respectively. Therefore, the index $\mathcal{A}_p(l)$, $p = 1, 2$ is a standardized distance between $\text{AV}(t)$ and $\text{ME}(t)$ according to ℓ^p Wasserstein distance for $p = 1, 2$.

Since this index is defined for each elementary level $[w_{l-1}, w_l]$, $\mathcal{A}_p(l)$ returns a set of level-wise skewness indices. This permits to evaluate positive or negative skewness at different levels of the set of histogram data. Furthermore, if $\mathcal{A}_p(l)$ is always positive (negative), for every l level, we propose a Total Skewness Index for a histogram variable expressed as follows:

$$\mathcal{A}_p = \sum_l \mathcal{A}_p(l).$$

5 Some Results on a Real Data Application

To illustrate the proposed statistics for histogram variables, we present an application on a histogram variable extracted from a real dataset. We have considered the values of the *water maximum temperatures* recorded by 21 meteorological Italian stations from 1, January, 2009 to 31, December 2009.³ For each station, the hourly data are summarized by a histogram, using an equi-width strategy. The 21 histograms are considered as the realization of a histogram variable. Finally, from each histogram we obtained the respective 21 quantile functions illustrated in Fig. 3a. It is possible to note that the temperature recorded in the period of observation ranges from 5 to 27 °C.

³The data are available at <http://www.mareografico.it>—Rete Mareografica Nazionale.

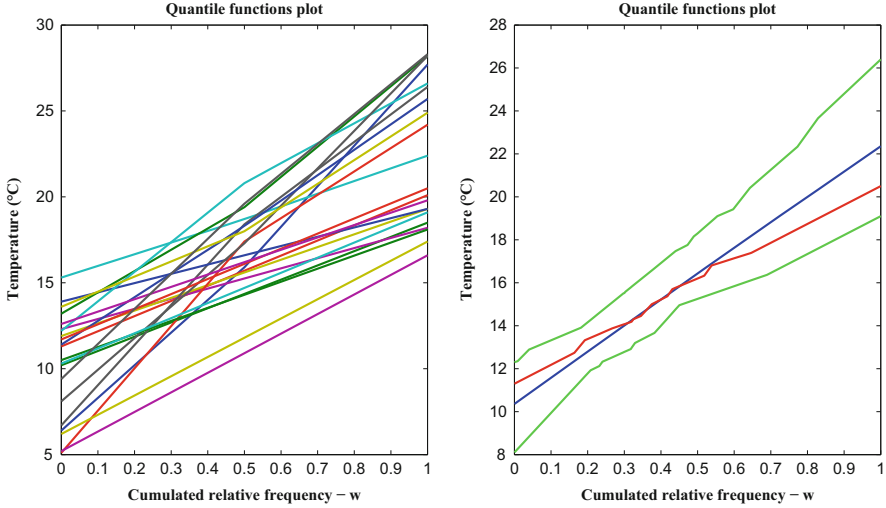


Fig. 3 (a) In the *left panel*, the plot of the quantile functions associated with histogram data of each meteorological station. (b) In the *right panel*, the median quantile function (in red), the first and the third quartile quantile functions (in green), and average quantile function (in blue)

Comparing the plot of the 21 quantile functions, the water temperatures recorded by some meteorological stations present a higher variability due to the different behavior between the coldest and the warmest periods.

By using the proposed strategy, we have computed for the set of 21 quantile functions and the $Q_{(\frac{1}{4}, N)}(t)$, $ME(t)$ and $Q_{(\frac{3}{4}, N)}(t)$ for the set of $N = 21$ histograms, and we have represented them in Fig. 3b. In particular, the median quantile function is obtained joining the selected $F_{(1l)}^{-1}(t)$ (with $w_{l-1} \leq t < w_l$) segments, where $l = 1, \dots, 37$.

The average and the median quantile functions intersect in some points as shown in Fig. 3b. In this case, we cannot consider a total skewness index but we can only compute the values $\mathcal{A}_p(l)$ for each interval of level $[w_l, w_{l+1})$. In particular, it is possible to observe that the median quantile function assumes higher values than the average quantile function ones in the interval $[0; 0.27)$ (showing a negative skewness), then after some oscillations in $[0.27; 0.34)$, the average quantile function always exceeds the median quantile function (showing a positive skewness).

In Fig. 4, it is evident the scale of the skewness by means of the level-wise skewness index (13) computed for $p = 1$ (the left panel) and $p = 2$ (the right panel). In both cases, the sign of $\mathcal{A}_p(l)$ highlights the different behavior of median and average quantile functions in the interval $[0, 1]$. Regarding their values, $\mathcal{A}_1(l)$ ranges between -0.12 and 0.45 , while $\mathcal{A}_2(l)$ ranges between -0.11 and 0.67 . In both cases, the skewness is stronger on the extremes than in the central part of the quantile functions. This result is consistent with the theory of quantile functions as reported also by Gilchrist [4].

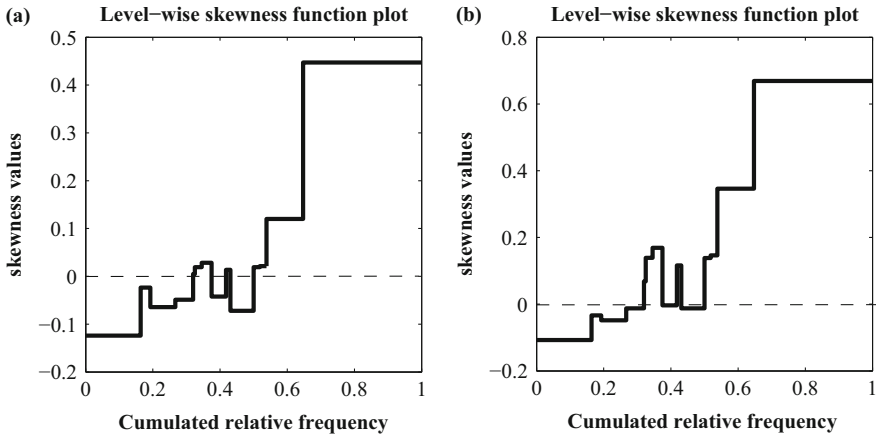


Fig. 4 (a) Level-wise skewness function for $p = 1$; (b) level-wise skewness function for $p = 2$

6 Conclusions

In this paper, we have presented a method for computing order statistics: the median, the first and the third quartile histogram for a histogram variable that takes a finite number of operations. In particular, we proposed to use a minimization of a distance criterion that is based on the quantile functions associated with the histogram data describing the units.

The order statistics here proposed are consistent with a level-wise order relationship. Consequently, the order-statistic quantile functions computed by our procedure may not correspond to an observed quantile function. However, order-statistics quantile functions are composed by pieces of quantiles functions associated with the observed histogram data. The proposed method does not guarantee a full order among quantile functions, but it is computationally affordable. This can be an advantage with respect to searching for a full order among a set of functions. Indeed, in the last case further restrictions and hypothesis are necessary on the set of the quantile functions. We have also proposed two measures of skewness for a histogram variable based on the ℓ_1 and the ℓ_2 Wasserstein distance. The application on the real dataset shows that the proposed order-statistics quantile functions, summarizing a set of data described by a histogram variable, is a useful exploratory tool and that the proposed skewness indices can detect sensible differences among the distributions of the quantiles associated with the histogram data.

References

1. Arroyo, J., Maté, C.: Forecasting histogram time series with k-nearest neighbours methods. *Int. J. Forecast.* **25**(1), 192–207 (2009)
2. Billard, L., Diday, E.: From the statistics of data to the statistics of knowledge: symbolic data analysis. *J. Am. Stat. Assoc.* **98**, 470–487 (2003)
3. Billard, L., Diday, E.: *Symbolic Data Analysis: Conceptual Statistics and Data Mining*. Wiley, Chichester (2006)
4. Gilchrist, W.: *Statistical Modelling with Quantile Functions*. Chapman and Hall, New York (2000)
5. González-Rivera, G., Arroyo, J.: Time series modelling of histogram-valued data: the daily histogram time series of S&P500 intradaily returns. *Int. J. Forecast.* **28**, 20–33 (2012)
6. Irpino, A., Verde, R.: Dynamic clustering of histograms using Wasserstein metric. In: Rizzi, A., Vichi, M. (eds.) *Advances in Computational Statistics*, pp. 869–876. Physica-Verlag, Heidelberg (2006)
7. Rüschendorf, L.: Wasserstein metric. In: Hazewinkel, M. (ed.) *Encyclopedia of Mathematics*. Springer, Berlin (2001)
8. Verde, R., Irpino, A.: Dynamic clustering of histograms using Wasserstein metric. In: Brito, P., et al. (eds.) *Selected Contributions in Data Analysis and Classification*, pp. 123–134. Springer, Berlin/Heidelberg (2007)
9. Verde, R., Irpino, A.: Comparing histogram data using a Mahalanobis-Wasserstein distance. In: Brito P. (ed.) *COMPSTAT 2008: Proceedings in Computational Statistics*, pp. 77–89. Physica-Verlag, Heidelberg (2008)

Rates for Bayesian Estimation of Location-Scale Mixtures of Super-Smooth Densities

Catia Scricciolo

Abstract

We consider Bayesian nonparametric density estimation with a Dirichlet process kernel mixture as a prior on the class of Lebesgue univariate densities, the emphasis being on the achievability of the error rate $n^{-1/2}$, up to a logarithmic factor, depending upon the kernel. We derive rates of convergence for the Bayes' estimator of *super-smooth* densities that are location-scale mixtures of densities whose Fourier transforms have sub-exponential tails. We show that a nearly parametric rate is attainable in the L^1 -norm, under weak assumptions on the tail decay of the true mixing distribution and the overall Dirichlet process base measure.

1 Introduction

Consider the estimation of a density f_0 on \mathbb{R} from observations X_1, \dots, X_n taking a Bayesian nonparametric approach. A prior is defined on a metric space of probability measures with Lebesgue density and a summary of the posterior, e.g., the posterior expected density, is employed. The so-called what if approach, which consists in investigating frequentist asymptotic properties of the posterior, under the non-Bayesian assumption that the data are generated from a *fixed* density, provides a way to validate priors on infinite-dimensional spaces. Desirable asymptotic properties of posterior distributions are consistency, minimax-optimal concentration rate of the posterior mass around the “truth” as the sample size grows, possibly with full adaptation to the regularity level of f_0 , if unknown, and distributional convergence.

C. Scricciolo (✉)
Bocconi University, Via Roentgen n. 1, 20136 Milan, Italy
e-mail: catia.scricciolo@unibocconi.it

For bounded and convex distances, posterior contraction rates yield upper bounds on convergence rates of the Bayes' estimator, thus motivating the interest in their study. Since the seminal articles of Ferguson [2] and Lo [4], the idea of constructing priors on spaces of densities by convoluting a fixed kernel with a random distribution has been successfully exploited in density estimation. Even if much progress has been done during the last decade in understanding frequentist asymptotic properties of mixture models, the choice of the kernel is a topic largely ignored in the literature, except for the article of Wu and Ghosal [9], mainly focussed on consistency. Posterior contraction rates for Dirichlet process kernel mixture priors have been investigated by Ghosal and van der Vaart [3] and Scricciolo [5]. One key message is that some constraints on the regularity of the kernel and on the tail decay of the true mixing distribution are necessary to accurately estimate a density. Most of the literature has dealt with the estimation of mixtures, with normal (or generalized normal) kernel and mixing distribution having either compact support or sub-exponential tails, finding a *nearly parametric* rate, up to a logarithmic factor, in the L^1 -distance, but there are almost no results beyond the Gaussian kernel. The aim of this work is to contribute to the understanding of the role of the kernel choice in density estimation with a Dirichlet process mixture prior. The main result states that a nearly parametric rate can be attained to estimate mixtures of super-smooth densities having Fourier transforms that decay exponentially, whatever the kernel tail decay, heavy tailed distributions, like Student's- t or Cauchy, being included, which have been proved to be extremely useful in accurately modeling different kinds of financial data. For example, individual stock indices can be modeled as stable laws. Multivariate stable laws have been fruitfully used in computer networks, see Bickson and Guestrin [1]. The assumption on the exponential tail decay of the true mixing distribution seems unavoidable in order to find a finite approximating mixture with a sufficiently restricted number of points. This step is a delicate mathematical point in the proof, see Lemma 1. Such an approximation result, which is reported in the Appendix, may be of autonomous interest as well. In Sect. 2, we fix the notation and present the result.

2 Main Result

We derive rates for location-scale mixtures of super-smooth densities. The model is $f_{F,G}(x) := \int_0^\infty (F * K_\sigma)(x) dG(\sigma)$, $x \in \mathbb{R}$, where K is a kernel density, $F \sim D_\alpha$ is a Dirichlet process with base measure $\alpha := \alpha(\mathbb{R})\bar{\alpha}$, for $0 < \alpha(\mathbb{R}) < \infty$ and $\bar{\alpha}$ a probability measure on \mathbb{R} , and $G \sim D_\beta$, with finite and positive base measure β on $(0, \infty)$. We assume that $f_0 = f_{F_0, G_0}$, with F_0 and G_0 denoting the true mixing distributions for the location and scale parameters, respectively. We use the following assumptions.

(A) The true mixing distribution G_0 for the scale parameter satisfies

$$\int_0^\infty \sigma \, dG_0(\sigma) < \infty \quad \text{and} \quad \int_0^\infty \frac{1}{\sigma} \, dG_0(\sigma) < \infty. \quad (1)$$

Also, for constants $d_1, d_2 > 0$ and $0 < \gamma_1^0, \gamma_2^0 \leq \infty$,

$$G_0(s) \lesssim e^{-d_1 s^{-\gamma_1^0}} \quad \text{as } s \rightarrow 0 \quad \text{and} \quad 1 - G_0(s) \lesssim e^{-d_2 s^{\gamma_2^0}} \quad \text{as } s \rightarrow \infty.$$

(B) The base measure β of the Dirichlet process prior for G has a continuous and positive Lebesgue density β' on $(0, \infty)$ such that, for constants $C_j, D_j > 0$, $j = 1, \dots, 4$, $q_1, q_2, r_1, r_2 \geq 0$ and $0 < \gamma_1, \gamma_2 \leq \infty$,

$$C_1 \sigma^{-q_1} e^{-C_2 \sigma^{-\gamma_1} (\log(1/\sigma))^{r_1}} \leq \beta'(\sigma) \leq C_3 \sigma^{-q_1} e^{-C_4 \sigma^{-\gamma_1} (\log(1/\sigma))^{r_1}} \quad (2)$$

for all σ in a neighborhood of 0, and

$$D_1 \sigma^{q_2} e^{-D_2 \sigma^{\gamma_2} (\log \sigma)^{r_2}} \leq \beta'(\sigma) \leq D_3 \sigma^{q_2} e^{-D_4 \sigma^{\gamma_2} (\log \sigma)^{r_2}} \quad (3)$$

for all σ large enough.

Remark 1 The right-hand side requirement in (1) has also been postulated by Tokdar [7], see condition 3 of Lemma 5.1 and condition 4 of Theorem 5.2, pp. 102–103. If, for example, G_0 is an $\text{IG}(\nu, \lambda)$, with shape parameter $\nu > 0$ and scale parameter $\lambda > 0$, then $\int_0^\infty \sigma^{-1} \, dG_0(\sigma) = (\nu/\lambda) < \infty$. If G_0 is a right-truncated distribution, then the requirement on the upper tail is satisfied with $\gamma_2^0 = \infty$. A right-truncated Inverse-Gamma distribution meets all the requirements of assumption (A).

Remark 2 Condition (2) is satisfied (with $r_1 = 0$) if β' is an Inverse-Gamma distribution. It can be seen that (2) implies that

$$\beta((0, s]) \leq \exp \left\{ -\frac{C_4}{2} s^{-\gamma_1} \left(\log \frac{1}{s} \right)^{r_1} \right\} \lesssim e^{-\frac{1}{2} C_4 s^{-\gamma_1}} \quad \text{as } s \rightarrow 0.$$

Condition (3) has been considered by van der Vaart and van Zanten [8], p. 2660, and implies that $\beta((s, \infty)) \lesssim \exp \{-D_4 s^{\gamma_2}/2\}$ as $s \rightarrow \infty$, see Lemma 4.9, p. 2669.

We assess rates for location-scale mixtures of symmetric stable laws. The result goes through to location-scale mixtures of Student's- t distributions.

Theorem 1 *Let K be the density of a symmetric stable law of index $0 < r \leq 2$. Suppose that $f_0 = \int_0^\infty (F_0 * K_\sigma) \, dG_0(\sigma)$, with the true mixing distribution F_0 for the*

location parameter satisfying the tail condition

$$F_0(\{\theta : |\theta| > t\}) \lesssim \exp\{-c_0 t^{1+I_{(1,2]}(r)/(r-1)}\} \quad \text{for large } t > 0, \quad (4)$$

for some constant $c_0 > 0$, and the true mixing distribution G_0 for the scale parameter satisfying assumption (A), with $\gamma_2^0 = \infty$. If the base measure α has a density α' such that, for constants $b > 0$ and $0 < \delta \leq 1 + I_{(1,2]}(r)/(r-1)$, satisfies

$$\alpha'(\theta) \propto e^{-b|\theta|^\delta}, \quad \theta \in \mathbb{R}, \quad (5)$$

the base measure β satisfies assumption (B), with $0 < \gamma_j \leq \gamma_j^0 \leq \infty$ and $\gamma_j < \gamma_j^0$ if $r_j > 0$, $j = 1, 2$, then the posterior rate of convergence relative to the Hellinger distance is $\varepsilon_n = n^{-1/2}(\log n)^\kappa$, with $\kappa > 0$ depending on γ_1^0 , γ_1 , γ_2 , and r .

Proof The proof is in the same spirit as that of Theorem 4.1 in Scricciolo [6], which, for space limitations, cannot be reported here. Let $\bar{\varepsilon}_n = n^{-1/2}(\log n)^\kappa$ and $\tilde{\varepsilon}_n = n^{-1/2}(\log n)^\tau$, with $\kappa > \tau > 0$ whose rather lengthy expressions we refrain from writing down. Let $0 < s_n \leq E(\log(1/\bar{\varepsilon}_n))^{-2\tau/\gamma_1}$, $0 < S_n \leq F(\log(1/\bar{\varepsilon}_n))^{2\tau/\gamma_2}$, and $0 < a_n \leq L(\log(1/\bar{\varepsilon}_n))^{2\tau/\delta}$, with $E, F, L > 0$ suitable constants. Replacing the expression of N in (A.19) of Lemma A.7 of Scricciolo [6], with that in Lemma 1, we can estimate the covering number of the sieve set

$$\mathcal{F}_n := \{f_{F,G} : F([-a_n, a_n]) \geq 1 - \bar{\varepsilon}_n/2, \quad G([s_n, S_n]) \geq 1 - \bar{\varepsilon}_n/2\}$$

and show that $\log D(\bar{\varepsilon}_n, \mathcal{F}_n, d_H) \lesssim (\log n)^{2\kappa} = n\bar{\varepsilon}_n^2$. Verification of the remaining mass condition $\pi(\mathcal{F}_n^c) \lesssim \exp\{-(c_2 + 4)n\bar{\varepsilon}_n^2\}$ can proceed as in the aforementioned theorem using, among others, the fact that $2\tau > 1$.

We now turn to consider the small ball probability condition. For $0 < \varepsilon < 1/4$, let $a_\varepsilon := (c_0^{-1} \log(1/(s_\varepsilon \varepsilon)))^{1/(1+I_{(1,2]}(r)/(r-1))}$ and $s_\varepsilon := (d_1^{-1} \log(1/\varepsilon))^{-1/\gamma_1^0}$. Let G_0^* be the re-normalized restriction of G_0 to $[s_\varepsilon, S_0]$, with S_0 the upper endpoint of the support of G_0 , and F_0^* the re-normalized restriction of F_0 to $[-a_\varepsilon, a_\varepsilon]$. Then, $\|f_{F_0^*, G_0^*} - f_0\|_1 \lesssim \varepsilon$. By Lemma 1, there exist discrete distributions $F'_0 := \sum_{j=1}^N p_j \delta_{\theta_j}$ on $[-a_\varepsilon, a_\varepsilon]$ and $G'_0 := \sum_{k=1}^N q_k \delta_{\sigma_k}$ on $[s_\varepsilon, S_0]$, with at most $N \lesssim (\log(1/\varepsilon))^{2\tau-1}$ support points, such that $\|f_{F'_0, G'_0} - f_{F_0^*, G_0^*}\|_\infty \lesssim \varepsilon$. For $T_\varepsilon := (2a_\varepsilon \vee \varepsilon^{-1/(r+I_{(0,1]}(r))))$,

$$\|f_{F'_0, G'_0} - f_{F_0^*, G_0^*}\|_1 \lesssim T_\varepsilon \|f_{F'_0, G'_0} - f_{F_0^*, G_0^*}\|_\infty + T_\varepsilon^{-r} \lesssim \varepsilon^{1-1/(r+I_{(0,1]}(r))}.$$

Without loss of generality, the θ_j 's and σ_k 's can be taken to be at least 2ε -separated. For any distribution F on \mathbb{R} and G on $(0, \infty)$ such that

$$\sum_{j=1}^N |F([\theta_j - \varepsilon, \theta_j + \varepsilon]) - p_j| \leq \varepsilon \quad \text{and} \quad \sum_{k=1}^N |G([\sigma_k - \varepsilon, \sigma_k + \varepsilon]) - q_k| \leq \varepsilon,$$

by the same arguments as in the proof of Theorem 4.1 in Scricciolo [6],

$$\|f_{F,G} - f_{F'_0, G'_0}\|_1 \lesssim \varepsilon.$$

Consequently,

$$\begin{aligned} d_{\text{H}}^2(f_{F,G}, f_0) &\leq \|f_{F,G} - f_{F'_0, G'_0}\|_1 + \|f_{F'_0, G'_0} - f_{F_0^*, G_0^*}\|_1 + \|f_{F_0^*, G_0^*} - f_0\|_1 \\ &\lesssim \varepsilon^{1-1/(r+I_{(0,1]}(r))}. \end{aligned}$$

By an analogue of the last part of the same proof, we get that $\pi(B_{\text{KL}}(f_0; \tilde{\varepsilon}_n^2)) \gtrsim \exp\{-c_2 n \tilde{\varepsilon}_n^2\}$.

Remark 3 Assumptions (4) on F_0 and (5) on α' imply that $\text{supp}(F_0) \subseteq \text{supp}(\alpha)$, thus, F_0 is in the *weak* support of D_α . Analogously, assumptions (A) on G_0 and (B) on β' , together with the restrictions on γ_j^0 , γ_j , $j = 1, 2$, imply that $\text{supp}(G_0) \subseteq \text{supp}(\beta)$, thus, G_0 is in the *weak* support of D_β .

Remark 4 If $\gamma_1 = \gamma_2 = \infty$, then also $\gamma_1^0 = \gamma_2^0 = \infty$, i.e., the true mixing distribution G_0 for σ is compactly supported on an interval $[s_0, S_0]$, for some $0 < s_0 \leq S_0 < \infty$, and (an upper bound on) the rate is given by $\varepsilon_n = n^{-1/2}(\log n)^\kappa$, with κ whose value for Gaussian mixtures ($r = 2$) reduces to the same found by Ghosal and van der Vaart [3] in Theorem 6.1, p. 1255.

Appendix

The following lemma provides an upper bound on the number of mixing components of finite location-scale mixtures of symmetric stable laws that uniformly approximate densities of the same type with compactly supported mixing distributions. We use \mathbb{E} and \mathbb{E}' to denote expectations corresponding to priors G and G' for the scale parameter Σ , respectively.

Lemma 1 *Let K be a density with Fourier transform such that, for constants A , $\rho > 0$ and $0 < r < 2$, $\Phi_K(t) = Ae^{-\rho|t|^r}$, $t \in \mathbb{R}$. Let $0 < \epsilon < 1$, $0 < a < \infty$ and $0 < s \leq S < \infty$ be given, with $(a/s) \geq 1$. For any pair of probability measures F on $[-a, a]$ and G on $[s, S]$, there exist discrete probability measures F' on $[-a, a]$ and G' on $[s, S]$, with at most*

$$N \lesssim \begin{cases} \frac{a}{s} \times \left(\frac{S}{s}\right)^r \left(\log \frac{1}{s\epsilon}\right)^{1+1/r}, & \text{if } 0 < r \leq 1, \\ \max \left\{ \left(\frac{a}{s}\right)^{r/(r-1)}, \left(\frac{S}{s}\right)^{r/(r-1)} \left(\log \frac{1}{s\epsilon}\right)^{1/(r-1)} \right\}, & \text{if } 1 < r < 2, \end{cases}$$

*support points, such that $\|\mathbb{E}[F * K_\Sigma] - \mathbb{E}'[F' * K_\Sigma]\|_\infty \lesssim \epsilon$.*

Proof We consider first the case where $1 < r < 2$ because, since $(a/s) \geq 1$ by assumption, we can appeal to Lemma A.1 of Scricciolo [6]. The arguments of the first part of the proof can be then used to deal also with the case where $0 < r \leq 1$. For each $s \leq \sigma \leq S$, since $\int_{-\infty}^{\infty} |\Phi_K(\sigma t)| dt < \infty$, the inversion formula can be applied to recover both $F * K_\sigma$ and $F' * K_\sigma$. For any $M > 0$ and $x \in \mathbb{R}$,

$$\begin{aligned} & |\mathbb{E}[(F * K_\Sigma)(x)] - \mathbb{E}'[(F' * K_\Sigma)(x)]| \\ &= \frac{1}{2\pi} \left| \int_s^S \int_{-\infty}^{\infty} e^{-itx} \Phi_K(\sigma t) \Phi_F(t) dt dG(\sigma) \right. \\ &\quad \left. - \int_s^S \int_{-\infty}^{\infty} e^{-itx} \Phi_K(\sigma t) \Phi_{F'}(t) dt dG'(\sigma) \right| \\ &= \frac{1}{2\pi} \left| \left(\int_{|t| \leq M} + \int_{|t| > M} \right) e^{-itx} [\Phi_F(t) \mathbb{E}[\Phi_K(\Sigma t)] - \Phi_{F'}(t) \mathbb{E}'[\Phi_K(\Sigma t)]] dt \right|. \end{aligned}$$

Let

$$U := \frac{1}{2\pi} \left| \int_{|t| \leq M} e^{-itx} [\Phi_F(t) \mathbb{E}[\Phi_K(\Sigma t)] - \Phi_{F'}(t) \mathbb{E}'[\Phi_K(\Sigma t)]] dt \right|$$

and

$$V := \frac{1}{2\pi} \left| \int_{|t| > M} e^{-itx} [\Phi_F(t) \mathbb{E}[\Phi_K(\Sigma t)] - \Phi_{F'}(t) \mathbb{E}'[\Phi_K(\Sigma t)]] dt \right|.$$

For $M \geq (\rho^{1/r} s)^{-1} (\log(1/(s^r \varepsilon)))^{1/r}$,

$$V \leq \frac{A}{2\pi} \int_{|t| > M} \int_s^S e^{-\rho(\sigma|t|)^r} d(G + G')(\sigma) dt \lesssim \varepsilon.$$

In order to find an upper bound on U , we apply Lemma A.1 of Ghosal and van der Vaart [3], p. 1260, to both F and G . There exists a discrete probability measure F' on $[-a, a]$, with at most $N_1 + 1$ support points, where N_1 is a positive integer to be suitably chosen later on, such that it matches the (finite) moments of F up to the order N_1 , i.e., $\mathbb{E}'[\Theta^j] = \mathbb{E}[\Theta^j]$ for all $j = 1, \dots, N_1$. Analogously, there exists a discrete probability measure G' on $[s, S]$, with at most N_2 support points, where N_2 is a positive integer to be suitably chosen later on, such that

$$\mathbb{E}'[\Sigma^{r\ell}] := \int_s^S \sigma^{r\ell} dG'(\sigma) = \int_s^S \sigma^{r\ell} dG(\sigma) =: \mathbb{E}[\Sigma^{r\ell}], \quad \ell = 1, \dots, N_2 - 1.$$

Both N_1 and N_2 will be chosen to be increasing functions of ε . In virtue of the latter matching conditions,

$$\begin{aligned} |\mathbb{E}[\Phi_K(\Sigma t)] - \mathbb{E}'[\Phi_K(\Sigma t)]| &\leq \left| \mathbb{E} \left[\Phi_K(\Sigma t) - A \sum_{\ell=0}^{N_2-1} \frac{(\rho(\Sigma|t|)^r)^\ell}{\ell!} \right] \right| \\ &\quad + \left| \mathbb{E}' \left[\Phi_K(\Sigma t) - A \sum_{\ell=0}^{N_2-1} \frac{(\rho(\Sigma|t|)^r)^\ell}{\ell!} \right] \right| \\ &\leq \frac{2A}{(N_2)!} (\rho(S|t|)^r)^{N_2}, \quad t \in \mathbb{R}. \end{aligned} \quad (6)$$

Using arguments of Lemma A.1 in Scricciolo [6] and inequality (6),

$$\begin{aligned} U &\leq \frac{1}{2\pi} \int_{|t| \leq M} |\Phi_F(t) \mathbb{E}[\Phi_K(\Sigma t)] - \Phi_{F'}(t) \mathbb{E}'[\Phi_K(\Sigma t)]| dt \\ &\leq \frac{1}{2\pi} \int_{|t| \leq M} \left| \Phi_F(t) - \sum_{j=0}^{N_1} \frac{(it)^j}{j!} \mathbb{E}[\Theta^j] \right| |\mathbb{E}[\Phi_K(\Sigma t)]| dt \\ &\quad + \frac{1}{2\pi} \int_{|t| \leq M} \left| \Phi_{F'}(t) - \sum_{j=0}^{N_1} \frac{(it)^j}{j!} \mathbb{E}'[\Theta^j] \right| |\mathbb{E}'[\Phi_K(\Sigma t)]| dt \\ &\quad + \frac{1}{2\pi} \sum_{j=0}^{N_1} \frac{|\mathbb{E}[\Theta^j]|}{j!} \int_{|t| \leq M} |t^j| |\mathbb{E}[\Phi_K(\Sigma t)] - \mathbb{E}'[\Phi_K(\Sigma t)]| dt \\ &\leq \frac{4Aa^{N_1}}{\pi r (\rho s^r)^{(N_1+1)/r}} \frac{\Gamma((N_1+1)/r)}{\Gamma(N_1+1)} + \frac{2A}{\pi (N_2)!} (1+aM)^{N_1} (\rho(SM)^r)^{N_2} M \\ &\sim \frac{4Aa^{N_1}}{\pi r (\rho s^r)^{(N_1+1)/r}} \frac{\Gamma((N_1+1)/r)}{\Gamma(N_1+1)} + \frac{\sqrt{2}A(1+aM)^{N_1} (\rho(SM)^r)^{N_2} M}{\pi^{3/2} e^{-N_2} N_2^{N_2-1/2}}, \end{aligned}$$

where, in the last line, we have used Stirling's approximation for $(N_2)!$, assuming N_2 is large enough. For $N_1 \lesssim \max\{\log(1/(s\varepsilon)), (a/s)^{r/(r-1)}\}$,

$$U_1 := \frac{4Aa^{N_1}}{\pi r (\rho s^r)^{(N_1+1)/r}} \frac{\Gamma((N_1+1)/r)}{\Gamma(N_1+1)} \lesssim \varepsilon.$$

Let M be such that $aM \geq 1$ and $(\rho^{1/r} SM) \geq 2a$. Then, for $N_2 \geq \max\{(2N_1+1)(r-1)/(r(2-r)), e^3(\rho^{1/r} SM)^{r/(r-1)}, \log(1/\varepsilon)\}$,

$$(1+aM)^{N_1} (\rho(SM)^r)^{N_2} M \leq (\rho^{1/r} SM)^{rN_2+2N_1+1} \leq (\rho^{1/r} SM)^{rN_2/(r-1)}$$

and

$$U_2 := \frac{\sqrt{2}A(1 + aM)^{N_1}(\rho(SM)^r)^{N_2}M}{\pi^{3/2}e^{-N_2}N_2^{N_2-1/2}} \lesssim \varepsilon.$$

Hence, $N_2 \lesssim \max\{(a/s)^{r/(r-1)}, ((S/s)^r)^{r/(r-1)} (\log(1/s^r\varepsilon))^{1/(r-1)}\}$.

In the case where $0 < r \leq 1$, since $(a/s) \geq 1$, we need to restrict the support of the mixing distribution F . To the aim, we consider a partition of $[-a, a]$ into $k = \lceil (a/s)(\log(1/(s\varepsilon)))^{1/r-1} \rceil$ subintervals I_1, \dots, I_k of equal length $0 < l \leq 2s(\log(1/(s\varepsilon)))^{-(1-r)/r}$ and, possibly, a final interval I_{k+1} of length $0 \leq l_{k+1} < l$. Let J be the number of intervals in the partition, which can be either k or $k + 1$. Write $F = \sum_{j=1}^J F(I_j)F_j$, where F_j denotes the re-normalized restriction of F to I_j . Then, for each $s \leq \sigma \leq S$, we have $(F * K_\sigma)(x) = \sum_{j=1}^J F(I_j)(F_j * K_\sigma)(x)$, $x \in \mathbb{R}$. For any probability measure F' such that $F'(I_j) = F(I_j)$, $j = 1, \dots, J$,

$$\begin{aligned} & |\mathbb{E}[(F * K_\Sigma)(x)] - \mathbb{E}'[(F' * K_\Sigma)(x)]| \\ & \leq \sum_{j=1}^J F(I_j) |\mathbb{E}[(F_j * K_\Sigma)(x)] - \mathbb{E}'[(F_j' * K_\Sigma)(x)]|, \quad x \in \mathbb{R}. \end{aligned}$$

Reasoning as in the case where $1 < r < 2$, with a to be understood as $l/2$ and N_1 as the number of support points of the generic F_j , for $M \geq ((\rho/2)^{1/r}s)^{-1}(\log(1/\varepsilon))^{1/r}$,

$$|\mathbb{E}[(F_j * K_\Sigma)(x)] - \mathbb{E}'[(F_j' * K_\Sigma)(x)]| \lesssim U + V \lesssim (U_1 + U_2) + \varepsilon, \quad x \in \mathbb{R}.$$

Since $(a/s) \lesssim (\log(1/(s\varepsilon)))^{-(1-r)/r}$ by construction, for $N_1 = \log(1/(s\varepsilon))$, it turns out that $U_1 \lesssim \varepsilon$. For $N_2 \geq \max\{N_1, 2e^4\rho(SM)^r \log(1/(s\varepsilon)), \log(1/\varepsilon)\}$,

$$(1 + aM)^{N_1}(\rho(SM)^r)^{N_2}M \leq M(2\rho(SM)^r \log(1/(s\varepsilon)))^{N_2}$$

and $U_2 \lesssim \varepsilon$. Then, $N_2 \lesssim (S/s)^r(\log(1/(s\varepsilon)))^2$ and the total number N^T of support points of F' is bounded above by

$$J \times N_1 \lesssim J \times N_2 \lesssim \frac{a}{s} \times \left(\frac{S}{s}\right)^r \left(\log \frac{1}{s\varepsilon}\right)^{1+1/r}.$$

The proof is thus complete.

Remark 5 Lemma 1 does not cover the case where $r = 2$, i.e., the kernel is Gaussian: this might possibly be due to the arguments laid out in the proof. This case can be retrieved from Lemma A.2 in Scricciolo [6] when $p = 2$.

References

1. Bickson, D., Guestrin, C.: Inference with multivariate heavy-tails in linear models. In: Lafferty, J.D., Williams, C.K.I., Shawe-Taylor, J., Zemel, R.S., Culotta, A. (eds.) *Advances in Neural Information Processing Systems*, vol. 23, pp. 208–216. Curran Associates, Inc., Red Hook (2010). <http://papers.nips.cc/paper/3949-inference-with-multivariate-heavy-tails-in-linear-models.pdf>
2. Ferguson, T.S.: Bayesian density estimation by mixtures of normal distributions. In: Rizvi, M.H., Rustagi, J.S., Siegmund, D. (eds.) *Recent Advances in Statistics*, pp. 287–302. Academic, New York (1983)
3. Ghosal, S., van der Vaart, A.W.: Entropies and rates of convergence for maximum likelihood and Bayes estimation for mixtures of normal densities. *Ann. Stat.* **29**, 1233–1263 (2001)
4. Lo, A.Y.: On a class of Bayesian nonparametric estimates: I. Density estimates. *Ann. Stat.* **12**, 351–357 (1984)
5. Scricciolo, C.: Posterior rates of convergence for Dirichlet mixtures of exponential power densities. *Electron. J. Stat.* **5**, 270–308 (2011)
6. Scricciolo, C.: Rates of convergence for Bayesian density estimation with Dirichlet process mixtures of super-smooth kernels. Working Paper No.1. DEC, Bocconi University (2011)
7. Tokdar, S.T.: Posterior consistency of Dirichlet location-scale mixture of normals in density estimation and regression. *Sankhyā* **68**, 90–110 (2006)
8. van der Vaart, A.W., van Zanten, J.H.: Adaptive Bayesian estimation using a Gaussian random field with inverse Gamma bandwidth. *Ann. Stat.* **37**, 2655–2675 (2009)
9. Wu, Y., Ghosal, S.: Kullback Leibler property of kernel mixture priors in Bayesian density estimation. *Electron. J. Stat.* **2**, 298–331 (2008)

Part II

Data Mining and Multivariate Data Analysis

Unsupervised Classification of Multivariate Time Series Data for the Identification of Sea Regimes

Mauro Bencivenga, Francesco Lagona, Antonello Maruotti,
Gabriele Nardone, and Marco Picone

Abstract

Unsupervised classification of marine data is helpful to identify relevant sea regimes, i.e. specific shapes that the distribution of wind and wave data takes under latent environmental conditions. We cluster multivariate marine data by estimating a multivariate hidden Markov model that integrates multivariate von Mises and normal densities. Taking this approach, we obtain a classification that accounts for the mixed (linear and circular) support of the observations, the temporal autocorrelation of the data, and the occurrence of missing values.

1 Introduction

Marine and atmospheric multivariate mixed data are disseminated by environmental protection agencies for a variety of purposes, ranging from the computation of simple descriptive summaries that communicate marine conditions to the public, to the estimation of sophisticated statistical models that detect air–sea interactions. Air–sea interaction models are exploited in several application areas, including studies of the drift of floating objects and oil spills, the design of off-shore structures, and studies of sediment transport and coastal erosion. These applications are especially important in coastal areas and semi-enclosed basins, where wind–

M. Bencivenga

ISPRA - Department for the Protection of Inland and Marine Waters, Via V. Brancati 60, 00144 Rome, Italy

F. Lagona (✉) • A. Maruotti

Department of Political Sciences, University Roma Tre, Via G. Chiabrera 199, 00145 Rome, Italy
e-mail: francesco.lagona@uniroma3.it

G. Nardone • M. Picone

ISPRA - Marine Service, Via V. Brancati 60, 00144 Rome, Italy

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,
Studies in Theoretical and Applied Statistics,

DOI 10.1007/978-3-319-27274-0_6

wave interactions are influenced by the orography of the site. Studies of air–sea interactions involve the analysis of multivariate, often incomplete, mixed-type time series of marine data, such as wind and wave measurement. These data are traditionally examined through numerical wind–wave models. Although well suited for oceans, numerical wind–wave models may provide inaccurate results under complex orography conditions. This has motivated the use of statistical models for the analysis of time series of wind and wave data [10].

We specify a multivariate hidden Markov model (MHMM) by describing wind–wave data in terms of latent environmental regimes, i.e., specific distributions that the data take under latent environmental conditions. This approach is particularly convenient under complex marine conditions, such as closed basins or coastal areas, where the correlation structure of the data can be decomposed according to a finite number of easily interpretable distributions. The literature on MHMM-based classification studies is dominated by Gaussian MHMMs for multivariate continuous data. MHMMs for data observed on different supports are less developed and traditionally specified by approximating the joint distribution of the data by a mixture whose components are products of univariate densities [7]. We extend this strand of the literature in the context of mixed linear and circular time series data. We focus in particular on pentavariate time series of wave and wind directions, wind speed, and wave height and period, typically collected to describe sea conditions in terms of wind–wave regimes. Previous work in this area includes mixtures whose components are specified as products of univariate or bivariate densities [4,5], which ignore temporal correlation, or hidden Markov models where single measurements are assumed as conditionally independent given the latent states of a Markov chain, which ignore correlation within latent classes [3]. We propose an MHMM where wind and wave directions are segmented by toroidal clusters, while (log-transformed) trivariate observations of wind speed and wave height and period are clustered within hyper-ellipses. Specifically, we approximate the joint distribution of the data by a product of bivariate von Mises densities [8] and trivariate normal densities, whose parameters evolve in time according to the states visited by a latent Markov chain [1, 6]. Under this setting, the transition probabilities matrix of the Markov chain captures regime-switching in time, accounting for temporal autocorrelation. This allows to cluster mixed linear and circular data separately, avoiding the definition of possibly hardly interpretable hyper-cylindrical clusters.

2 Wind–Wave Data in Adriatic Sea

Our analysis is based on semi-hourly environmental profiles with three linear and two circular components: wind speed, significant wave height, and wave period, wind direction and wave direction. The data were recorded in wintertime, in the period 18/1/2011–9/3/2011 by the buoy of Ancona, located in the Adriatic Sea at about 30 km from the coast. According to NOAA definitions, the significant wave height is the average crest-to-trough height of the highest one-third waves in a wave record and it be estimated from the variance of a wave elevation record assuming that the nondirectional spectrum is narrow. The wave period is the average time

between corresponding points on a wave profile passing a measurement location. Wave period is strictly correlated with wave height. Under conditions of calm sea, wave period depends only on the wave speed. As the height increases, waves start moving faster, and the period increases. The relationship between wind speed and wave period is weak and depends on the wind direction.

To account for the cumulative effect that wind has on waves, wind data were smoothed by taking, for each measurement, the average of wind speeds and the circular average of wind directions, observed during the last 6 h. Buoy data often include missing values because of transmission errors or malfunctioning of the device. About 5 % of the profiles includes at least one missing value. We assume that missing values occur at random. Under this hypothesis, the contribution of missing patterns to the likelihood can be ignored, facilitating model-based clustering of the data.

In the literature, such as wave atlas, marine data are typically depicted in terms of univariate distributions. The complex wind–wave interaction structure in the Adriatic Sea is, however, better shown by Fig. 1, which displays the scatter plots of

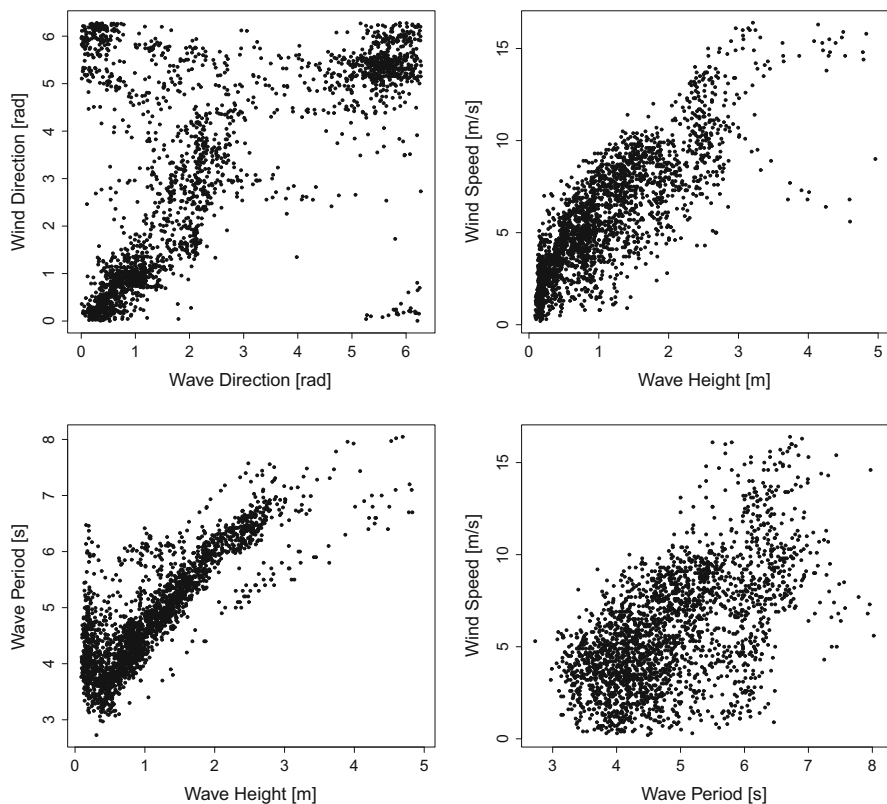


Fig. 1 Wind and wave direction (0 represents North) as well as wind speed, wave height and period, observed at a buoy of the Adriatic Sea in wintertime

the circular and the linear available observations. For simplicity, bivariate circular data are plotted on the plane, although data points are actually on a torus. Point coordinates are measured in radian, with North as the origin ($0 - 2\pi$). Although a number of patterns appear in these scatter plots, their interpretation is difficult due to the weak correlation of the circular measurements and the skewness of the linear observations, traditionally explained as the result of the complex orography of the Adriatic Sea and often held responsible for the inaccuracy of numerical wind-wave models. Nevertheless, the observations might result from the mixing of a number of latent environmental regimes, conditionally on which the distribution of the data takes a shape that is easier to interpret than the shape taken by the marginal distributions. By taking a HMM approach, we cluster directional and planar data separately to account for the different nature of the data, and simultaneously pair these clusters into a number of latent classes evolving in time according to a Markov chain, being interpretable as time-varying regimes of air-sea interactions.

3 A Gaussian-von Mises MHMM for Linear-Circular Data

The data considered in this paper are in the form of a time series $\mathbf{z}_{0:T} = (\mathbf{z}_t, t = 0, \dots, T)$, with $\mathbf{z}_t = (\mathbf{x}_t, \mathbf{y}_t)$ where $\mathbf{x}_t = (x_{1t}, x_{2t}) \in [0, 2\pi)^2$ are the bivariate circular components and $\mathbf{y}_t = (y_{1t}, y_{2t}, y_{3t}) \in \mathbb{R}^3$ are log-transformed linear components. In HMM-based classification studies, the temporal evolution of class membership is driven by a latent Markov chain, which can be conveniently described as a multinomial process in discrete time. Accordingly, we introduce a sequence $\boldsymbol{\xi}_{0:T} = (\boldsymbol{\xi}_t, t = 0, \dots, T)$ of multinomial variables $\boldsymbol{\xi}_t = (\xi_{t1} \dots \xi_{tK})$ with one trial and K classes, whose binary components represent class membership at time t . The joint distribution $p(\boldsymbol{\xi}_{0:T}; \boldsymbol{\pi})$ of the chain is fully known up to a parameter $\boldsymbol{\pi}$ that includes K initial probabilities $\pi_k = P(\xi_{0k} = 1)$, $k = 1, \dots, K$, $\sum_k \pi_k = 1$, and K^2 transition probabilities $\pi_{hk} = P(\xi_{tk} = 1 | \xi_{t-1,h} = 1)$, $h, k = 1, \dots, K$, $\sum_k \pi_{hk} = 1$. Formally, we assume that

$$p(\boldsymbol{\xi}_{0:T}; \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{\xi_{0k}} \prod_{t=1}^T \prod_{h=1}^K \prod_{k=1}^K \pi_{hk}^{\xi_{t-1,h} \xi_{tk}}. \quad (1)$$

The specification of a multivariate HMM is completed by assuming that the observations are conditionally independent, given a realization of the Markov chain. As a result, the conditional distribution of the observed process, given the latent process, takes the form of a product density, say

$$f(\mathbf{z}_{0:T} | \boldsymbol{\xi}_{0:T}) = \prod_{t=0}^T \prod_{k=1}^K (f_k(\mathbf{z}_t))^{\xi_{tk}},$$

where $f_k(\mathbf{z})$, $k = 1, \dots, K$ are K multivariate densities. For classification purposes, these densities are usually assumed to be known up to a number of parameters that indicate the locations and the shapes of K clusters. We assume that circular and linear observations are conditionally independent given a realization of the Markov chain, and introduce a family of bivariate densities $f(\mathbf{x}; \boldsymbol{\beta})$ on the torus, indexed by a parameter $\boldsymbol{\beta}$, and a family of trivariate densities on the plane, $f(\mathbf{y}; \boldsymbol{\gamma})$, indexed by a parameter $\boldsymbol{\gamma}$. Formally, we assume that

$$f(\mathbf{z}_{0:T} | \boldsymbol{\xi}_{0:T}) = \prod_{t=0}^T \prod_{k=1}^K (f(\mathbf{x}_t | \boldsymbol{\beta}_k) f(\mathbf{y}_t | \boldsymbol{\gamma}_k))^{\xi_{tk}}. \quad (2)$$

Integrating $f(\mathbf{z}_{0:T} | \boldsymbol{\xi}_{0:T}) p(\boldsymbol{\xi}_{0:T})$ with respect to $\boldsymbol{\xi}_{0:T}$, we obtain the marginal distribution of the observed data, known up to a parameter $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\beta}, \boldsymbol{\gamma})$, on which our classification procedure is based. In particular, we first maximize the likelihood function

$$L(\boldsymbol{\theta}; \mathbf{z}_{0:T}) = \sum_{\boldsymbol{\xi}_{0:T}} p(\boldsymbol{\xi}_{0:T}; \boldsymbol{\pi}) f(\mathbf{z}_{0:T} | \boldsymbol{\xi}_{0:T}), \quad (3)$$

and find the maximum likelihood estimate $\hat{\boldsymbol{\theta}}$. Secondly, we cluster the data according to the posterior probabilities of class membership

$$\hat{p}_{ik} = P(\xi_{ik} = 1 | \mathbf{z}_{0:T}; \hat{\boldsymbol{\theta}}) = \mathbb{E}(\xi_{ik} | \mathbf{z}_{0:T}; \hat{\boldsymbol{\theta}}), \quad (4)$$

based on $\hat{\boldsymbol{\theta}}$. The computational complexity of the estimation of both $\hat{\boldsymbol{\theta}}$ and \hat{p}_{ik} depends on the choice of the multivariate densities which are used to model the circular and the linear components of the observations. We exploit bivariate von Mises and trivariate normal densities, described below, as a compromise between numerical complexity and modeling flexibility.

The bivariate von Mises density in the form introduced by Mardia et al. [8] is a parametric distribution on the torus, which naturally embeds the bivariate normal distribution when the range of observations is small. Its density is given by

$$f(\mathbf{x}; \boldsymbol{\beta}) = \frac{\exp(\beta_{11} \cos(x_1 - \beta_1) + \beta_{22} \cos(x_2 - \beta_2) + \beta_{12} \sin(x_1 - \beta_1) \sin(x_2 - \beta_2))}{C(\boldsymbol{\beta})}, \quad (5)$$

with normalizing constant

$$C(\boldsymbol{\beta}) = 4\pi^2 \sum_{m=0}^{\infty} \binom{2m}{m} \left(\frac{\beta_{12}^2}{4\beta_{11}\beta_{22}} \right)^m I_m(\beta_{11}) I_m(\beta_{22}),$$

where $I_m(x)$ is the modified Bessel function of order m .

The univariate marginal densities $f(x_i; \boldsymbol{\beta})$ $i = 1, 2$ depend on the marginal mean angles β_i and on shape parameters β_{12} , β_{ii} , and β_{ij} with $i, j = 1, 2$ and $i \neq j$. For $\beta_{12} = 0$, x_1 and x_2 are independent and each of them follows a von Mises distribution with marginal mean angles β_i and marginal concentrations β_{ii} .

The conditional distributions $f(x_i|x_j; \boldsymbol{\beta})$ $i, j = 1, 2$ $i \neq j$ are von Mises with parameters depending on $\boldsymbol{\beta}$.

To model the joint distribution of (log-transformed) wave height and period and wind speed, we use a family of K trivariate normal densities

$$f(\mathbf{y}; \boldsymbol{\gamma}) = N \left(\begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \end{pmatrix}, \begin{pmatrix} \gamma_{11} & \gamma_{12} & \gamma_{13} \\ \gamma_{12} & \gamma_{22} & \gamma_{23} \\ \gamma_{13} & \gamma_{23} & \gamma_{33} \end{pmatrix} \right). \quad (6)$$

4 Likelihood Inference

As our data are in the form of incomplete profiles, we refer to $\mathbf{z}_{t,\text{mis}} = (\mathbf{x}_{t,\text{mis}}, \mathbf{y}_{t,\text{mis}})$ and $\mathbf{z}_{t,\text{obs}} = (\mathbf{x}_{t,\text{obs}}, \mathbf{y}_{t,\text{obs}})$, respectively, indicate the missing and observed parts of a circular-linear observation at time t .

If the data are missing at random (MAR), the missing data mechanism can be ignored and the maximum likelihood estimate of parameter $\boldsymbol{\theta}$ is the maximum point of the marginal likelihood function

$$L(\boldsymbol{\theta} | \mathbf{z}_{0:T,\text{obs}}) = \sum_{\boldsymbol{\xi}_{0:T}} p(\boldsymbol{\xi}_{0:T}; \boldsymbol{\pi}) \prod_{t=0}^T \int f(\mathbf{z}_t | \boldsymbol{\xi}_t; \boldsymbol{\beta}, \boldsymbol{\gamma}) d\mathbf{z}_{t,\text{mis}} \quad (7)$$

$$f(\mathbf{z}_t | \boldsymbol{\xi}_t; \boldsymbol{\beta}, \boldsymbol{\gamma}) = f(\mathbf{x}_t | \boldsymbol{\xi}_t; \boldsymbol{\beta}) f(\mathbf{y}_t | \boldsymbol{\xi}_t; \boldsymbol{\gamma}).$$

In order to estimate $\boldsymbol{\theta}$, we maximize $L(\boldsymbol{\theta})$ by using a version of the EM algorithm. EM algorithms are based on the definition of a complete-data log-likelihood function, obtained by considering the sampling distribution of both the observed and the unobserved quantities. As our HMM is a mixture which integrates circular and normal densities, the unobserved quantities are not only the missing measurements, but also the unknown class memberships. As a result, the complete-data log-likelihood function can be defined as follows

$$\begin{aligned} \log L_{\text{comp}}(\boldsymbol{\theta}, \boldsymbol{\xi}_{0:T}, \mathbf{z}_{0:T}) &= \sum_{k=1}^K \xi_{0k} \log \pi_k + \sum_{t=1}^T \sum_{h=1}^K \sum_{k=1}^K \xi_{t-1,h} \xi_{t,k} \log \pi_{hk} \\ &+ \sum_{t=0}^T \sum_{k=1}^K \xi_{tk} \log f(\mathbf{x}_t; \boldsymbol{\beta}_k) + \sum_{t=0}^T \sum_{k=1}^K \xi_{tk} \log f(\mathbf{y}_t; \boldsymbol{\gamma}_k). \end{aligned} \quad (8)$$

The algorithm is iterated by alternating the expectation (E) and maximization (M) step. Given the estimate $\hat{\theta}_s$, obtained at the end of the s -th iteration, the $(s+1)$ -th iteration is initialized by an E-step, which evaluates the expected value of (8) with respect to the conditional distribution of the missing values given the observed data. For the HMM, this distribution takes a complex, but tractable, form, because it factorizes as follows:

$$f(\xi_{0:T}, z_{0:T, \text{mis}} | z_{0:T, \text{obs}}; \hat{\theta}_s) = p(\xi_{0:T} | z_{0:T, \text{obs}}; \hat{\theta}_s) f(z_{0:T, \text{mis}} | \xi_{0:T}, z_{0:T, \text{obs}}; \hat{\theta}_s) \quad (9)$$

where

$$f(z_{0:T, \text{mis}} | \xi_{0:T}, z_{0:T, \text{obs}}; \hat{\theta}_s) = f(\mathbf{x}_{0:T, \text{mis}} | \xi_{0:T}, \mathbf{x}_{0:T, \text{obs}}; \hat{\beta}_s) \\ \times f(\mathbf{y}_{0:T, \text{mis}} | \xi_{0:T}, \mathbf{y}_{0:T, \text{obs}}; \hat{\gamma}_s)$$

and

$$f(\mathbf{x}_{0:T, \text{mis}} | \xi_{0:T}, \mathbf{x}_{0:T, \text{obs}}; \hat{\beta}_s) = \prod_{t=0}^T \prod_{k=1}^K \left(f(\mathbf{x}_{t, \text{mis}} | \xi_{tk} = 1, \mathbf{x}_{t, \text{obs}}; \hat{\beta}_{ks}) \right)^{\xi_{tk}} \\ f(\mathbf{y}_{0:T, \text{mis}} | \xi_{0:T}, \mathbf{y}_{0:T, \text{obs}}; \hat{\gamma}_s) = \prod_{t=0}^T \prod_{k=1}^K \left(f(\mathbf{y}_{t, \text{mis}} | \xi_{tk} = 1, \mathbf{y}_{t, \text{obs}}; \hat{\gamma}_{ks}) \right)^{\xi_{tk}}. \quad (10)$$

Each distributions $f(\mathbf{x}_{t, \text{mis}} | \xi_{tk} = 1, \mathbf{x}_{t, \text{obs}}; \hat{\beta}_{ks})$ and $f(\mathbf{y}_{t, \text{mis}} | \xi_{tk} = 1, \mathbf{y}_{t, \text{obs}}; \hat{\gamma}_{ks})$ in (10) are equal to 1 if the observed profile of circular or linear variables at time t is complete; it is otherwise equal to the multivariate circular or linear distribution, evaluated at $\beta = \hat{\beta}_{ks}$ or $\gamma = \hat{\gamma}_{ks}$, if both circular or linear measurements at time t are missing; it finally reduces to the circular or linear conditional distributions, evaluated at $\beta = \hat{\beta}_{ks}$ or $\gamma = \hat{\gamma}_{ks}$, if one observation is missing.

The factorization (9) facilitates the evaluation of the expected complete-data log-likelihood, which can be computed in terms of iterated expectations as follows

$$Q(\theta | \hat{\theta}_s) = \mathbb{E} \left(\log L_{\text{comp}}(\theta, \xi_{0:T}, z_{0:T} | z_{0:T, \text{obs}}; \hat{\theta}_s) \right) \\ = \sum_{k=1}^K \mathbb{E}(\xi_{0k} | z_{0:T, \text{obs}}, \hat{\theta}_s) \log \pi_k \quad (11)$$

$$+ \sum_{t=1}^T \sum_{h=1}^K \sum_{k=1}^K \mathbb{E}(\xi_{t-1, h} \xi_{tk} | z_{0:T, \text{obs}}, \hat{\theta}_s) \log \pi_{h,k} \quad (12)$$

$$+ \sum_{t=0}^T \sum_{k=1}^K \mathbb{E}(\xi_{tk} | \mathbf{z}_{0:T,\text{obs}}, \hat{\boldsymbol{\theta}}_s) \mathbb{E}(\log f(\mathbf{x}_t; \boldsymbol{\beta}_k) | \mathbf{x}_{t,\text{obs}}, \hat{\boldsymbol{\beta}}_{ks}) \quad (13)$$

$$+ \sum_{t=0}^T \sum_{k=1}^K \mathbb{E}(\xi_{tk} | \mathbf{z}_{0:T,\text{obs}}, \hat{\boldsymbol{\theta}}_s) \mathbb{E}(\log f(\mathbf{y}_t; \boldsymbol{\gamma}_k) | \mathbf{y}_{t,\text{obs}}, \hat{\boldsymbol{\gamma}}_{ks}). \quad (14)$$

The M-step of the algorithm updates the estimate $\hat{\boldsymbol{\theta}}_s$ with a new estimate $\hat{\boldsymbol{\theta}}_{s+1}$, which maximizes the above function Q . This function is the sum of three functions that depend on independent sets of parameters and can thus be then maximized separately. Maximization of (12) with respect to the transition probabilities π_{hk} provides the closed-form updating formula

$$\hat{\pi}_{hk(s+1)} = \frac{\sum_{t=1}^T \hat{p}_{t-1,t,hk}(\hat{\boldsymbol{\theta}}_s)}{\sum_{t=1}^T \hat{p}_{t-1,h}(\hat{\boldsymbol{\theta}}_s)}, \quad h, k = 1, \dots, K.$$

While the maximization of (14) reduces to a battery of exact updating equations [2], maximization of (13) reduces to K separate nonlinear systems of five equations, which may be solved following, e.g., the iterative procedure suggested by Mardia et al. [9].

5 Results

We clustered the data described in Sect. 2 by fitting a Gaussian–von Mises HMMs with $K = 3$ states, chosen according to the Bayesian information criterion. Table 1 displays the maximum likelihood estimates of the model. Figure 2 shows the contours of the toroidal densities and the trivariate densities, represented as three bivariate densities, for the three latent regimes. The data points are allocated according to the most probable regime. The model detects three regimes of straightforward interpretation.

The first component of the model is associated with Sirocco and weak Bora episodes. Wind and wave directions appear strictly synchronized. Waves travel driven by winds that blow from a similar directional angle, so when wind blows from the south, wave travels southeasterly along the major axis of the basin. When the wind turns along the north direction, under cyclonic atmospheric circulation, wave comes from northeast. As a result, wave can reach significant height, winds contribute to increase the velocity of waves, so the mean period is high.

A similar phenomenon is captured by the second component of the model, although it is originated by an anticyclonic atmospheric circulation. In this regime, northern Bora jets generate high waves that travel along the major axis of the basin. Compared to the other two regimes, waves and winds are highly concentrated around one modal direction. Most of the wind energy is transferred to the sea surface

Table 1 Estimated parameters of a three-state multivariate hidden Markov model with mixed (log) linear and circular components

			State 1	State 2	State 3
Circular parameters	Wave mean direction	β_1	1.335	0.150	0.727
	Wind mean direction	β_2	1.312	0.060	4.466
	Wave directional concentration	β_{11}	4.101	8.673	0.024
	Wind directional concentration	β_{22}	1.230	6.601	1.629
	Wave/wind directional dependence	β_{12}	2.812	9.519	-1.535
Linear parameters	Wave average height	γ_1	0.019	0.372	-1.306
	Wave average period	γ_2	1.600	1.641	1.400
	Wind average speed	γ_3	1.445	2.044	0.985
	Wave height variance	γ_{11}	0.253	0.228	0.261
	Wave period variance	γ_{22}	0.037	0.032	0.020
	Wind speed variance	γ_{33}	0.161	0.084	0.299
	Height/period covariance	γ_{12}	0.083	0.080	-0.014
	Height/speed covariance	γ_{13}	0.134	0.120	0.221
	Period/speed covariance	γ_{23}	0.034	0.039	-0.013
	Destination state		1	2	3
	Origin state	1	0.969	0.026	0.005
		2	0.012	0.985	0.003
		3	0.005	0.004	0.991
	Initial state		0.000	0.000	1.000

and, as a result, the correlation between wind speed and wave height and period is larger than that observed under Sirocco or Maestral episodes. As expected, most of the profiles with the highest waves in the sample are clustered under this regime.

The third component is associated with periods of calm sea: weak winds generate small waves with low periods. Under this regime, the shape of the joint distribution of wave and wind directions is essentially spherical and relate to northwesterly Maestral episodes. As expected, wind and wave directions are poorly synchronized under this regime, because wave direction is more influenced by marine currents than by wind direction during weak wind episodes.

The rows at the bottom of Table 1 include the estimated transition probabilities and initial probabilities of the latent Markov chain. As expected, the transition probability matrix is essentially diagonal, reflecting the temporal persistence of the states.

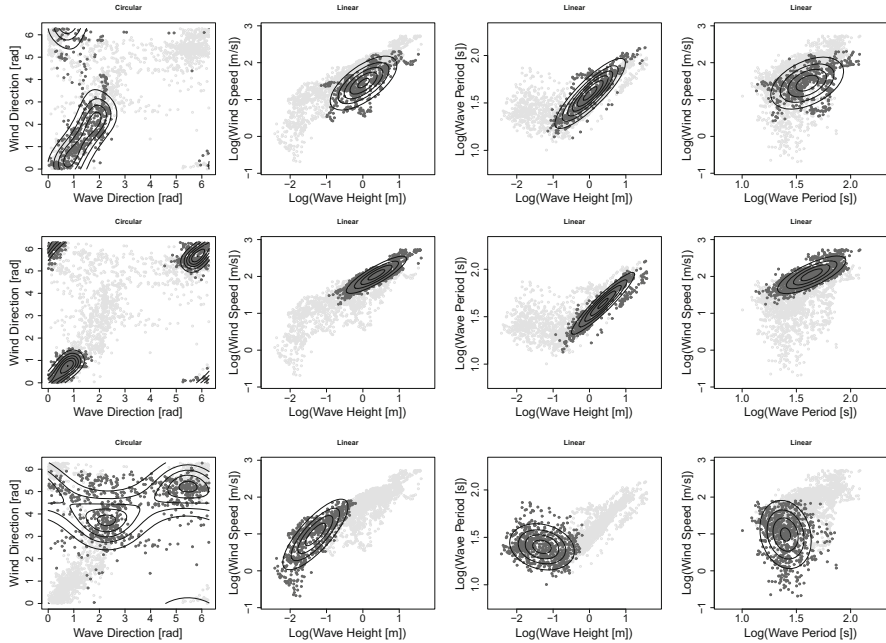


Fig. 2 Circular and linear components of a three-state hidden Markov models

6 Discussion

We illustrate a HMM-based classification method for multivariate mixed-type time series, focusing on linear and circular variables. The data are clustered according to trivariate normal and bivariate von Mises distributions, which are associated to the states of a latent Markov chain. Our classification procedure is motivated by issues that arise in marine studies, but can be easily adapted to a wide range of real-world cases, including, for example, ecological studies of animal behavior, where direction and speed of movements are recorded [3], and bioinformatics applications, where sequences of protein dihedral angles [8] are recorded with a number of continuous variables.

The model is useful to cluster marine data in a number of latent classes or regimes, associated with toroidal and elliptical clusters. Classification is carried out by accounting for both the temporal autocorrelation of the data and the special structure of the circular data. The combination of multivariate von Mises and normal distributions allows for a simple specification of the dependence structure between variables and for the computational feasibility of a mixture-based classification strategy where missing values can be efficiently handled within a likelihood framework.

References

1. Bulla, J., Lagona, F., Maruotti, A., Picone, M.: A multivariate hidden Markov model for the identification of sea regimes from incomplete skewed and circular time series. *J. Agric. Biol. Environ. Stat.* **17**, 544–567 (2012)
2. Cappé, O., Moulines, E., Rydén, T.: *Inference in Hidden Markov Models*. Springer, New York (2005)
3. Holzmann, H., Munk, A., Suster, M., Zucchini, W.: Hidden Markov models for circular and linear-circular time series. *Environ. Ecol. Stat.* **13**, 325–347 (2006)
4. Lagona, F., Picone, M.: A latent-class model for clustering incomplete linear and circular data in marine studies. *J. Data Sci.* **9**, 585–605 (2011)
5. Lagona, F., Picone, M.: Model-based clustering of multivariate skew data with circular components and missing values. *J. Appl. Stat.* **39**, 927–945 (2012)
6. Lagona, F., Picone, M.: Maximum likelihood estimation of bivariate circular hidden Markov models from incomplete data. *J. Stat. Comput. Simul.* **83**, 1223–1237 (2013)
7. Lagona, F., Maruotti, A., Picone, M.: A non-homogeneous hidden Markov model for the analysis of multi-pollutant exceedances data. In: Dymarsky, P. (ed.) *Hidden Markov Models, Theory and Applications*, pp. 207–222. InTech, Rijeka (2011)
8. Mardia, K.V., Hughes, G., Taylor, C.C., Singh, H.: A multivariate von Mises distribution with applications to bioinformatics. *Can. J. Stat.* **36**, 99–109 (2007)
9. Mardia, K., Taylor, C., Subramaniam, G.: Protein bioinformatics and mixtures of bivariate von Mises distributions for angular data. *Biometrics* **63**, 505–512 (2007)
10. Monbet, V., Ailliot, P., Prevosto, M.: Survey of stochastic models for wind and sea-state time series. *Probab. Eng. Mech.* **22**, 113–126 (2007)

An Evaluation of the Student Satisfaction Based on CUB Models

Barbara Cafarelli and Corrado Crocetta

Abstract

In 2009 the Faculty of Economics of Foggia started a project called “Analisi della student satisfaction” with the aim of creating an internal quality control system based on students’ feedback. Every year a customer satisfaction survey is carried out, where all students attending lectures are asked to evaluate the services provided. This paper presents an evaluation of the student satisfaction over the last two academic years. In order to understand the level of satisfaction and the psychological mechanism behind the students’ evaluation process an approach based on CUB models is adopted. At the end, multidimensional scaling (MDS) techniques were used to investigate the existence of student subgroups with similar attitudes towards the Faculty’s services and the overall satisfaction, highlighting the eventually presence of similarities with the latent variables estimated by CUB models.

1 Introduction

In the last few decades the demand for information from policy makers and stakeholders has increased considerably. A large part of the information available is referred to data collected on objective basis but there is an increasing need for subjective data and in particular for customer satisfaction feedback.

In the university system there is strong interest from potential students, families and institutions for information about the quality and the reputation of the university.

B. Cafarelli (✉) • C. Crocetta

Department of Economics, University of Foggia, Foggia, Italy
e-mail: barbara.cafarelli@unifg.it; corrado.crocetta@unifg.it

This consideration led the Faculty of Economics of Foggia to develop the student satisfaction analysis project with the aim of building a monitoring system devoted to improving the quality of the services provided (see [1]).

Since the Academic Year 2009–2010, a customer survey has been carried out every Academic Year, paying particular attention to the respondents' perceptions and to the underlying psychological construct behind them.

In order to assess how students' judgement is influenced by personal feelings towards the items under investigation and by the inherent uncertainty associated with the choice of the ordinal values featuring on the questionnaire responses (see [3]), the CUB models (see [4, 5]) were used.

We hereby propose to detect significant similarities and differences between the raters' overall judgment by comparing the estimated CUB models, with the aid of a multidimensional scaling (MDS).

The MDS approach is also proposed in order to confirm the presence of the two latent CUB variables and the role they play in the satisfaction process.

In this paper we only report the results of the satisfaction survey referring to the services most frequently used by the students: Logistics (*L*), Registrar's Office (*RO*), Teaching (*T*), Website (*WS*) plus the overall satisfaction with the Faculty (*OS*) for two Academic Years (A.Y.): 2009–2010 and 2010–2011.

The paper is organized as follows: after a brief description of the survey's characteristics and the presentation of the methodology applied, the results will be discussed and scrutinized in detail.

2 The Student Survey

The evaluation of student satisfaction with respect to the services provided has been performed by means of a questionnaire proposed to all the undergraduates attending lectures at the Faculty of Economics of the University of Foggia during the A.Y. 2009–2010 and 2010–2011. We collected 968 interviews for the first year analysed and 832 for the second year. The respondents are the 35 and 30 % of the total of registered students, respectively. All interviewed students completed the questionnaire.

The services under evaluation were: Logistics, Registrar's Office, Teaching, Website, Laboratories, Library and Job Placements.

In this study, the evaluations of the services library and job placement are not considered because they are not directly managed by the Faculty of Economics.

For each service we considered different items and then included a question about the overall satisfaction level of the service. In particular, for each service considered the students were asked about tangibles, reliability, responsiveness, assurance and empathy.

The students gave a score to each service considered and to the overall satisfaction which was expressed in a Likert scale from 1 (extremely unsatisfied) to 7 (extremely satisfied).

The questionnaire also included questions about general characteristics of interviews, areas to improve and reasons for enrollment.

In order to detect potential difficulties of respondents the questionnaire was pre-tested.

3 The Statistical Procedure

The student satisfaction assessment was made by CUB models to understand how customers’ preferences are influenced by a subjective personal feeling towards the items under investigation and by the inherent uncertainty associated with the choice of the ordinal values featuring on the questionnaire responses.

In particular, D’Elia and Piccolo [4] suggest modelling feeling using a shifted binomial random variable and uncertainty by using a discrete uniform random variable. In this way, a CUB model (see [4, 5]) is obtained as a mixture distribution combining these two random variables as follows:

$$P(R = r) = \pi \binom{m-1}{r-1} (1-\xi)^{r-1} \xi^{m-r} + (1-\pi) \frac{1}{m} \quad r = 1, 2, \dots, m \quad (1)$$

where $\xi \in [0, 1]$, $\pi \in [0, 1]$, r is the rating with $m > 3$. For each service under judgment, sample data $(r_1, \dots, r_i, \dots, r_n)'$ were observed, where r_i is the rating expressed by the i th subject and n is the number of the observed subjects, and were considered as a realization of the random sample $(R_1, \dots, R_i, \dots, R_n)$.

The ξ parameter indicates the importance of a latent aspect that can be called feeling, whereas the π parameter measures the importance of uncertainty in the final judgement. The interpretation of these two parameters depends on the scores of judgements. In this case a Likert scale from 1 (extremely unsatisfied) to 7 (completely satisfied) is adopted. As a consequence, the ξ parameter is the inverse of feeling/agreement with the item. Low values correspond to high levels of feeling. The π parameter represents the uncertainty of the choice and high values correspond to low levels of uncertainty of the respondents.

The statistical procedure proposed consists of the following steps:

1. Model (1) was fitted for each considered service for the Academic Years 2009–2010 and 2010–2011. 10 CUB models were estimated. Then the profiles of satisfaction, estimated by CUB models, were compared with those observed.

For this data, traditional goodness-of-fit indices cannot be properly applied, because they detect significant differences even though there is an “almost perfect fit” (see [3]). For this reason, to assess the goodness of fit of each estimated model, the dissimilarity index was used (see [3, 6]):

$$Diss = \frac{1}{2} \sum_{r=1}^7 \left| f_r - p_r(\hat{\xi}, \hat{\pi}) \right| \quad (2)$$

where f_r are the observed relative frequencies and $p_r(\hat{\xi}, \hat{\pi})$ are the probabilities estimated by CUB model. The dissimilarity index ranges between 0 and 1 (where values less than 0.1 suggest a very good fit).

The models (1) were estimated using software CUB.R (3.0) implemented in R (see [6]).

The presence of similarity or dissimilarity between the judgements of each item was also investigated by comparing the representation in the parametric space of the estimated values of feeling and uncertainty. ξ and π are characterized by different variability patterns and play a different role in determining the shape of the estimated distributions. For this reason the use of the Euclidean distance might cause misleading interpretations of the CUB models in the parametric space (see [7]). To avoid this problem, the Kullback–Leibler (*KL*) divergence was used to evaluate dissimilarity among the estimated rating distributions (see [10]).

The Kullback–Leibler statistics (see [8]) is defined as follows:

$$KL = \frac{N_i N_j}{N_i + N_j} \left[\sum_x (p(x, \theta_i) - p(x, \theta_j)) \ln \frac{p(x, \theta_i)}{p(x, \theta_j)} \right]_{\theta_i = \hat{\theta}_i; \theta_j = \hat{\theta}_j} \quad (3)$$

where $p(x, \theta_i)$ is the probability distribution function, i and j are the items under comparison and N_i and N_j are, respectively, the number of observations for the items i and j . In this study i and j vary from 1 to the number of items considered for both Academic Years.

The *KL* statistics (3) was used to test the null hypothesis $H_0 : \theta_i = \theta_j$ against $H_1 : \theta_i \neq \theta_j$, where the vector parameters $\theta = (\xi, \pi)$ have been replaced by the maximum likelihood estimators (see [8]). In these cases, the *KL* divergence follows the χ^2 test with two degrees of freedom (see [9]);

2. A non-parametric MDS approach, based on *KL* divergences, was then applied to study the relevant dissimilarities between the services analyzed and to highlight the two latent variables estimated by CUB models (see [2]).
3. The results of MDS were compared with the procedure for clustering ordinal data proposed by Corduas [8]. In particular, the comparison was done with the hierarchical cluster analysis based on *KL* divergence with complete linkage method.

4 The Results of the Student Satisfaction Survey

Before presenting the results of the procedure described in Sect. 3, it can be useful to have a look at the observed data. The students interviewed expressed a good level of satisfaction with the services under investigation and the overall satisfaction for both considered years. In fact, mean, median and mode were equal or higher than 4.0 on a 7 point scale.

Table 1 Median, mode, mean and Fisher–Pearson coefficient of skewness for the overall satisfaction for the analysed services and for the Overall satisfaction for the academic years 2009–2010 and 2010–2011

Service		Academic years	
		2009–2010	2010–2011
Teaching	Mean	4.9	4.9
	Median	5.0	5.0
	Mode	5.0	5.0
	Skewness	−0.9	−1.0
Website	Mean	5.3	5.0
	Median	6.0	5.0
	Mode	6.0	5.0
	Skewness	−1.1	−0.9
Registrar’s office	Mean	4.3	4.0
	Median	5.0	4.0
	Mode	5.0	5.0
	Skewness	−0.5	−0.3
Logistics	Mean	4.3	3.9
	Median	4.0	4.0
	Mode	5.0	5.0
	Skewness	−0.5	−0.4
Overall satisfaction	Mean	4.8	4.6
	Median	5.0	5.0
	Mode	5.0	5.0
	Skewness	−0.9	−0.6

The observed distributions have a negative skewness as shown by the negative values of the Fisher–Pearson coefficients (Table 1).

The students were more satisfied with Teaching and Website than Registrar’s office and Logistics, the situation seems to get worse in 2010–2011 with respect to the previous year.

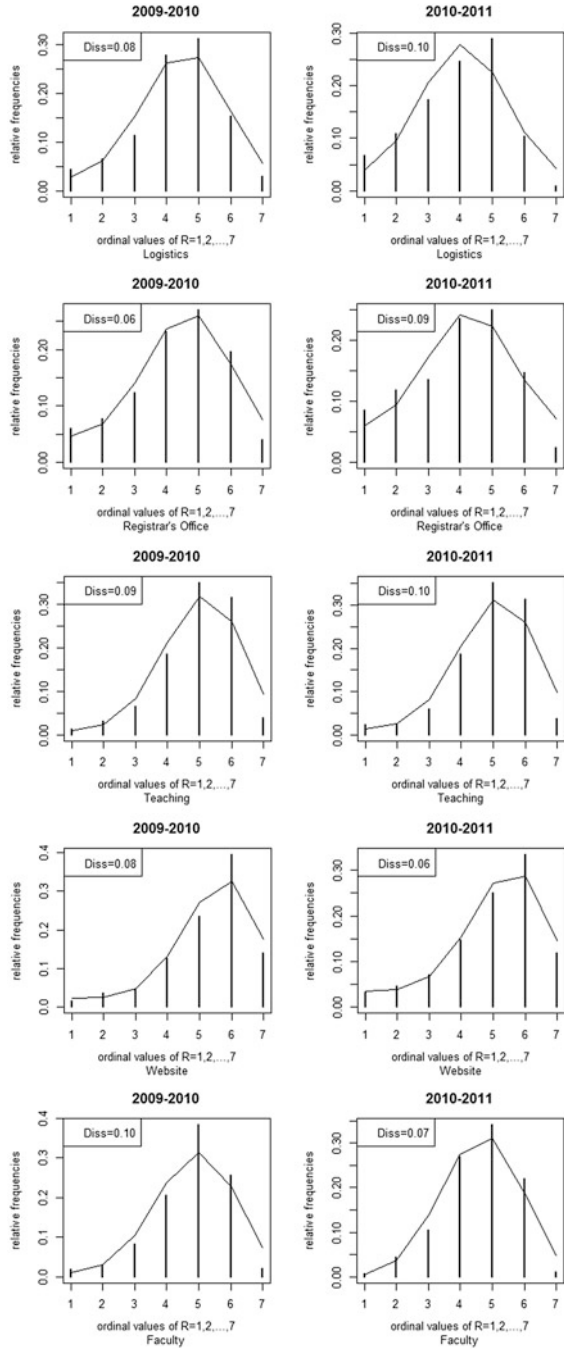
In order to investigate these results, paying attention to the respondents’ perceptions and the underlying psychological construct behind them, 10 CUB models were fitted.

The comparison between observed and estimated profiles of satisfaction showed that CUB models well-represent the observed distributions as also suggested by the values of the dissimilarity index of each estimated CUB model that were always less or equal to 0.1 (Fig. 1).

By dividing in four regions of the parametric space ($[0, 1] \times [0, 1]$) of $\theta = (\xi, \pi)$, four different areas are obtained denoting, respectively, low level of feeling with low level of uncertainty (1st quadrant in the upper right), low level of feeling with high level of uncertainty (2nd quadrant, in the upper left), high level of feeling with low level of uncertainty (3rd quadrant, in the lower left), high level of feeling with low level of uncertainty (4th quadrant, in the lower right).

Figure 2 shows that all the parameters estimated by CUB models are in the 4th quadrant confirming the good level of satisfaction and a low level of uncertainty among the respondents in both Academic Years. In particular, the feeling parameter

Fig. 1 Observed and estimated profiles of satisfaction and *Diss* values for logistics, registrar’s office, teaching, website and overall satisfaction for academic years 2009–2010 and 2010–2011



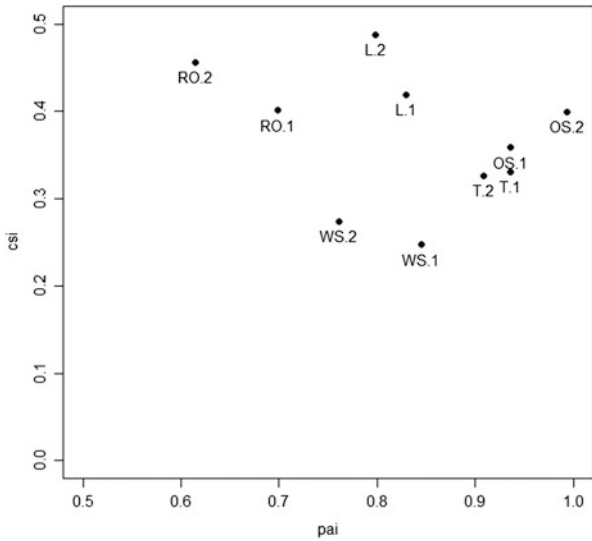


Fig. 2 Feeling and uncertainty estimated by CUB models for logistics (*L*), registrar’s office (*RO*), teaching (*T*), website (*WS*) plus the overall satisfaction with the faculty (*OS*) for academic year 2009–2010 denoted by 0.1 and academic year 2010–2011 denoted by 0.2

$\hat{\xi}$ ranges from 0.25 to 0.42 in the first year and between 0.27 and 0.49 in the second year. The estimated uncertainty parameters, $\hat{\pi}$, are greater than 0.69 for the first year and 0.61 for the second, which means that the interviewees are quite precise when it comes to giving marks. In order to better describe the results obtained we displayed only the fourth quadrant; in this way, the most appreciated services are the ones nearest to the lower right corner of Fig. 2.

If we consider the Academic Year 2009–2010, with respect to the feeling, the best services are Website, Teaching and Overall satisfaction for the Faculty. Registrar’s Office and Logistic have a lower level of feeling. The students are more confident about Overall satisfaction, Teaching and Website and a little less for Logistic and Registrar’s Office.

In the second year the examined services show worse performances in terms of feeling compared to the first year, with the only exception of Teaching, and an increasing amount of respondents’ uncertainty. It is interesting to point out that the overall satisfaction with respect to the Faculty services is very close to Teaching. That confirms the descriptive results, that quality of teaching strongly influences the perception of the respondents.

We can summarize that the most appreciated services are Website and Teaching whereas Logistic and Registrar’s Office have lower feeling.

The similarities between teaching and overall satisfaction and between Registrar’s Office and Logistic are confirmed by (3) (p -value > 0.05).

The differences between the performance of each service in the 2 years are also tested by (3). In the second year, the reduction of feeling and uncertainty for Logistics and Registrar’s Office is significant (p -value < 0.05). The decrease of feeling and the increase of uncertainty for the Overall satisfaction are significant (p -value < 0.05).

The differences between the 2 years with regard to Website’s and Teaching performances are not significant.

The graphical proximity between the estimated points in the parametric space must be checked by using an appropriate technique, in order to avoid the risk of misleading interpretation of data, arising by the use of the Euclidean distance (see [7, 9]).

For this reason, following Corduas [8], we analyzed similarities between profiles of satisfaction estimated by CUB models, using a hierarchical cluster method based on the Kullback–Leibler divergence, using the complete linkage approach (Fig. 3).

Four main groups were detected: Website for the two Academic Years, Registrar’s Office and Logistics for 2009–2010, Teaching for both Academic Years and overall satisfaction with the Faculty for A.Y. 2009–2010 and Logistic, Registrar’s Office for A.Y. 2009–2010 and the overall satisfaction for 2010–2011.

The use of the hierarchical cluster however, did not give information about the students’ perceptions in terms of feeling and uncertainty.

To this purpose (see [8]), a non-metric MDS approach based on KL divergence was used to assess the presence of dissimilarities between the performance of

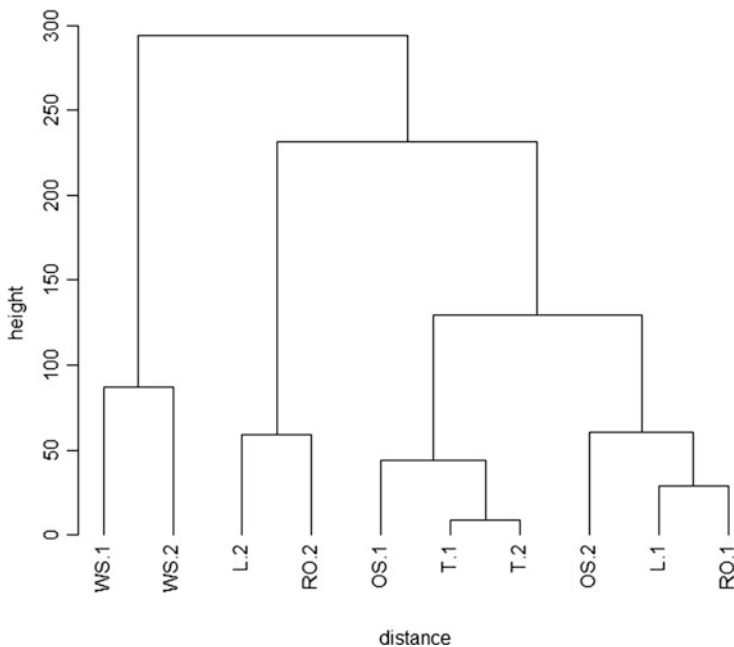


Fig. 3 Hierarchical cluster: complete linkage method

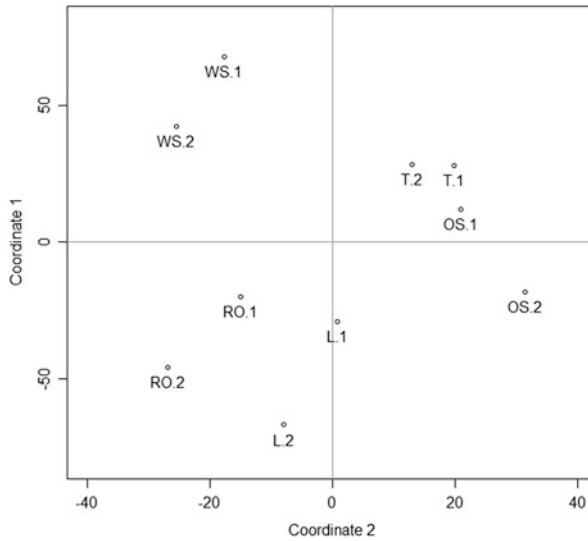


Fig. 4 MDS map

the various services provided and the overall quality of the Faculty over the two investigated years.

This approach allows us to detect relevant differences among the services issued, in terms of feeling and uncertainty, paying attention to the psychological mechanisms behind the evaluation process and to analyse the two latent aspects of CUB models (Fig. 4).

The results obtained confirm the ones obtained by the cluster analysis. The values of the stress index suggests considering only two dimensions.

The MDS approach shows the existence of student subgroups with similar attitudes in terms of feeling and uncertainty towards the Faculty’s services and the overall satisfaction and the presence of the two latent dimensions in the psychological construct that shapes the respondents’ behaviour, as suggested by the CUB based analysis.

In fact, the horizontal dimension (coordinate 1) clearly shows Website on one side and Logistics and Registrar’s Office on the other and could be seen as the feeling, whereas the vertical dimension (coordinate 2) shows Logistics, Registrar’s Office and Website on one hand and Teaching and overall satisfaction with the Faculty on the other and could be seen as the uncertainty level of respondents (Fig. 4).

It is easy to verify that Fig. 4 is similar to Fig. 2 rotated of 90°. We can conclude that the MDS analysis confirms results obtained by CUB models.

5 Concluding Remarks

In this study we focussed on an integrated approach as a tool to account for the psychological mechanisms behind the students' evaluation process and to compare ratings expressed by the students on different services.

In particular, CUB models were used to understand the role played by two unobservable components, supposed to be behind the evaluation process: feeling and uncertainty. The MDS was used to facilitate the interpretation of these latent components in the parametric space. This approach can represent an option to the procedure suggested by Corduas [8]. In fact, in this study we have seen that the conclusions achieved by applying the two integrated procedures (cluster analysis and the MDS on the CUB model results) are similar.

However, the use of the MDS has a further advantage. In fact, this technique led us to detect not only the presence of statistically similar groups among the respondents, but also to highlight the underlying dimensions that allow the researchers to explain observed similarities or dissimilarities between the investigated items. This procedure is also useful to compare distributions in terms of skewness, median and mode.

The results of the proposed method have been used by the Faculty of Economics of the University of Foggia to improve services and consequently student satisfaction levels. They confirm the central role of the quality of teaching with respect to the overall satisfaction.

Acknowledgments This paper is part of the MIUR project PRIN 2008 “Modelli per variabili latenti basati su dati ordinali: metodi statistici ed evidenze empiriche” (Research Unit of Napoli Federico II)—CUP E61J10000020001 and of the project “Analisi della student satisfaction”, University of Foggia (2009–2012). The authors jointly designed and realized the work here described. However, B. Cafarelli edited Sects. 2, 3 and 4 and C. Crocetta edited Sects. 1 and 5.

References

1. Cafarelli, B., Crocetta, C., Spada, A.: La student satisfaction relativa ai servizi della Facoltà di Economia dell'Università degli studi di Foggia. In: *Annali del Dipartimento di Scienze Statistiche “CARLO CECCHI”*, vol. IX, pp. 221–244, CLEUP, Padova (2010). ISBN: 978886129665
2. Cafarelli, B., Pilone, V., Conte, A., Gammariello, D., Del Nobile, M.A.: Development of consumer acceptable products using CUB analysis: an example with burgers from Dairy Cattle. *J. Sens. Stud.* 30(5), 413–424 (2015)
3. Iannario, M., Piccolo, D.: A new statistical model for the analysis of customer satisfaction. *Qual. Technol. Quant. Manag.* 7, 149–168 (2010)
4. D'Elia, A., Piccolo, D.: A mixture model for preference data analysis. *Comput. Stat. Data Anal.* 49, 917–934 (2005)
5. Iannario, M.: Modelling shelter choices in a class of mixture models for ordinal responses. *Stat. Methods Appl.* 21, 1–22 (2012)
6. Iannario, M., Piccolo, D.: A short guide CUB 3.0 Program. www.researchgate.net/publication/260952393_A_Short_Guide_to_CUB_3.0_Program (2014)

7. Corduas, M., Iannario, M., Piccolo, D.: A class of statistical models for evaluating services and performances. In: Bini, M., Monari, P., Piccolo, D., Salmaso, L. (eds.) *Statistical Methods for the Evaluation of Educational Services and Quality of Products*. Contributions to Statistics, pp. 99–117. Physica-Verlag, Heidelberg (2010)
8. Corduas, M.: A statistical procedure for clustering ordinal data. *Quaderni di Statistica* **10**, 177–187 (2008)
9. Piccolo, D.: On the moments of a mixture of uniform and shifted binomial random variables. *Quaderni di Statistica* **5**, 85–104 (2003)
10. Kullback, S.: *Information Theory and Statistics*. Wiley, New York (1959)

Dimensions of Well-Being and Their Statistical Measurements

Carla Ferrara, Francesca Martella, and Maurizio Vichi

Abstract

Nowadays, a relevant challenge regards the assessment of a global measure of well-being by using composite indicators of different features such as level of wealth, comfort, material goods, standard living, quality and availability of education, etc.

In this paper, we focus on statistical methodologies designed to build composite indicators of well-being by detecting latent components and assessing the statistical relationships among indicators. We will consider some constrained versions of Principal Component Analysis (PCA) which allow to specify disjoint classes of variables with an associated component of maximal variance. Once the latent components are detected, a Structural Equation Model (SEM) has been used to evaluate their relationships. These methodologies will be compared by using a data set from 34 member countries of the Organization for Economic Co-operation and Development (OECD).

1 Introduction

In recent years, the assessment of the subjective perception of overall well-being of citizens has received increasing attention. Governments and economists are aware that Gross Domestic Product (GDP) and other traditional measures of economic progress fail to measure the quality of the life. In fact, well-being depends on a number of factors not related exclusively to the economic and material elements but also to lifestyle, food choice, health condition, and environment.

C. Ferrara (✉) • F. Martella • M. Vichi

Department of Statistical Sciences, Sapienza University of Rome, Piazzale Aldo Moro, 5, 00185 Rome, Italy

e-mail: carla.ferrara@uniroma1.it; francesca.martella@uniroma1.it; maurizio.vichi@uniroma1.it

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,

Studies in Theoretical and Applied Statistics,

DOI 10.1007/978-3-319-27274-0_8

In this complex empirical framework, the methodological statistical problem is to find a unique measure, i.e., a *composite indicator* synthesizing various factors of well-being from social to economic, naturally arises. Several methodologies have been proposed to define composite indicators (Saisana et al. [1]). In their methodological literature, two aspects have been mainly investigated: (1) the identification of key indicators to be used; (2) the ways in which these indicators can be brought together to make a coherent system of information. Problems occur in how to choose, aggregate, and weight the single observed variables among an available suite, or in how to identify the components driving the composite indicators. Therefore, the selection of the weights and the way the variables are combined are very sensitive and important issues to be studied in order to exclude an oversimplification of a complex system and prevent potentially misleading signals.

Dimensional reduction methods, such as Principal Component Analysis (PCA) or Factor Analysis (FA), are the most natural tools to compute a composite indicator of a set of observed variables. It can be defined as a linear combination of the observed variables which explains the largest part of their total variance, i.e., it is the first principal component or factor, provided that a unique relevant principal component—in terms of explained variance (e.g., a unique component with variance greater or equal to 1)—is defined. However, frequently the observed variables specify more than one pertinent component, as it is for the global indicator of well-being, representing a complex latent variable (LV) synthesized by several other latent variables (such as living standards, health conditions, and environment). Thus, there is a need to synthesize the latent variables into a unique final composite indicator and to assess the relationships among latent and observed variables. A Structural Equation Model (SEM) can be used for these last two purposes. One assumption of SEM is that observed variables are partitioned in blocks of variables, each related to a well a priori identified latent concept. As Kline [2] suggests: a priori does not mean exclusively confirmatory, in the sense that in most of situations the researcher needs to test more than one model. Frequently, the model chosen by the researcher is not supported by the data and it must be rejected or the hypotheses on which it is based must be modified. Furthermore, when no a priori knowledges are available, it can be useful an Exploratory Factor Analysis (EFA) to discover LVs in order to identify a model to be tested. In this case, blocks of variables that identify latent components and explain the largest part of the total variance of the observed variables need to be computed. To reach this aim, variables are partitioned into blocks, generally, by means of PCA or FA. EFA, traditionally, has been used to explore the possible underlying factor structure of a set of observed variables without imposing a preconceived structure on the outcome. By performing EFA, the underlying factor structure is identified. On the other hand, Confirmatory Factor Analysis (CFA) developed by Jöreskog [3] is a statistical technique used to verify the factor structure of a set of observed variables. CFA allows the researcher to test the hypothesis that a relationship between observed variables and their underlying latent constructs exists. The researcher uses knowledge of the theory, empirical research, or both, postulates the relationship pattern a priori and then tests the hypothesis statistically.

However, there are some limitations of PCA and FA, in the correct assignment of an original variable to a given component, because frequently variables have a relevant loading for more than one component and therefore there is no clear classification of these variables into a single component. This problem has also relevant implications in the interpretation of components of PCA and FA and for this reason frequently a rotation, such as varimax, is introduced to simplify the structure of component loadings matrix. However, rotation frequently does not solve this classification problem.

The problem of rotating the loading matrix \mathbf{A} ($J \times Q$) into a particular meaningful rotation is linked to the concept of simple structure, introduced by Thurstone [4]. The general idea is that factors have real meaning when many variables do not depend on all of them. Thus, the loading matrix should have as many zero coefficients as possible to show a simplest structure and presumably to allow the most meaningful interpretation. Thurstone specifies the following properties: (1) each row of \mathbf{A} shall have at least one zero; (2) each column of \mathbf{A} shall have at least Q zeros, where Q is the number of LVs; (3) for every pair of columns of \mathbf{A} there shall be several rows in which one loading is zero and one is nonzero; (4) for every pair of columns of \mathbf{A} , a large proportion of rows shall have two zero loadings (if $Q \geq 4$); (5) for every pair of columns of \mathbf{A} , there shall preferably be only a small number of rows with two nonzero coefficients.

Reiersol [5] investigates and modifies these conditions in order to allow factor model identification. Let $\Sigma = \mathbf{A}\Phi\mathbf{A}' + \Psi$ be the resulting covariance matrix of a factor model, where Φ and Ψ are the covariance matrices of factors and errors, respectively. He assumes that there are at least Q zeros in each column of \mathbf{A} and that (1) the rank of \mathbf{A}_m is $Q - 1$, where \mathbf{A}_m is the matrix \mathbf{A} with zeros in the m^{th} column; (2) the rank of each submatrix obtained by deleting a row of \mathbf{A}_m is $Q - 1$, and (3) the rank of each matrix obtained by adding a row of \mathbf{A}_m , not contained in \mathbf{A}_m , is Q . If Ψ is identified, a necessary and sufficient condition for its identification is that \mathbf{A} does not contain any other submatrices satisfying (1), (2), and (3).

In an EFA, where no hypotheses concerning the factors are involved, it is convenient to choose these restrictions so that $\Phi = \mathbf{I}$ and $\mathbf{A}'\Psi^{-1}\mathbf{A}$ is diagonal. In a CFA, on the other hand, some relationships between factors and variables are hypothesized and therefore restrictions on certain elements of \mathbf{A} and Φ are required in advance. For example, $a_{jm} = 0$ means that the m^{th} factor does not load on the j^{th} variable. Thus, Jöreskog [3] defines unrestricted and restricted solutions, whether or not the common factor space is restricted. The concept of restricted model has been also analyzed by Zou et al. [6]. They propose a new method called sparse principal component analysis (SPCA), which shrinks elements of \mathbf{A} towards zero, by imposing the elastic net constraint on the loadings and, therefore, deriving sparse loadings.

In this paper we propose three methods to partition the variables, named: step-wise PCA, restricted PCA, and disjoint PCA. After blocks of variables are defined SEM has been applied. Two different approaches for estimating SEM parameters can be considered: the covariance-based and the component-based techniques. The first approach, which includes maximum likelihood estimation method (ML-SEM

or LISREL, [3]), has been for many years the only estimation method aiming at reproducing the sample covariance matrix of the observed variables through the model parameters. On the other hand, the second approach, also known as PLS Path Modeling (PLS-PM, [7]), was developed as an alternative approach to LISREL, as a more flexible technique for the treatment of a huge amount of data characterized by small sample sizes with respect to the number of variables, less restrictive assumptions are required as compared to classical covariance-based approaches in terms of distributions and measurement scales. It provides estimates of the LVs in such way that they are the most correlated with each other, according to a path diagram structure, and the most representative of each corresponding block of manifest variables (MVs). Among others, recently Tenenhaus and Tenenhaus [8] show that PLS-PM, except for some specific cases, maximizes a function of covariance among connected LVs and under two different type of constraints is put on outer weights and on latent variable scores. In this paper we will use this second approach.

In particular, the paper is organized as follows. Section 2 describes constrained PCA methods and a constrained version of clustering and disjoint PCA (CDPCA, [9]) model. In Sect. 3, Structural Equation Model (SEM, [10]) is introduced in order to assess relationships among variables and, in particular, to estimate a network of causal relationships linking two or more latent complex concepts, each measured through a number of observable indicators [11]. In Sect. 4, a motivating example, based on a real data set described in [12], is illustrated to show the implementation of the PLS-PM. In particular, we analyze 34 member countries considering well-being under three aspects: material living conditions, quality of life, sustainability, similarly to the paper proposed by [13]. Concluding remarks and discussion are given in the last section.

2 Constrained PCA Methods

Three different methodologies are described to partition the original variables into blocks with a component of maximal variance for each block.

2.1 Stepwise PCA

The first method, named “stepwise PCA”, consists of performing several iterated PCA steps on the relevant Q principal components (PCs) (assuming scaled variables, those with eigenvalue greater or equal to 1), and deleting at each PCA iteration the variable which takes the greatest loading on the last relevant PC. This procedure is repeated until a single component, which defines a subset of variables, is obtained. This strategy is iterated for the matrix containing all the variables excluded previously, until all the variables are associated to a single PC. In this way, disjoint

PCs are obtained, since each variable can be exclusively linked to only a PC. The algorithm is represented in the following steps:

- step 1 Compute the first Q PCs on Σ_X to obtain $\mathbf{Y} = \mathbf{X}\mathbf{A}^{(Q)}$; $\mathbf{B} = \mathbf{I}_Q$
- step 2 $m = \operatorname{argmax}\{|a_{lQ}^{(Q)}|: l = 1, \dots, J\}$, find the variable with the largest loading, in absolute value on the Q^{th} PC
- step 3 $\mathbf{B}(m, m) = 0$, delete the variable m
- step 4 Compute PCA on $\mathbf{B}\Sigma_X\mathbf{B}$ and repeat Steps 2 and 3 until \mathbf{Y} is a single PC linear combination of a subset of variables that represent the first cluster
- step 5 On the remaining variables repeat Steps 2, 3, and 4, until a partition of variables is obtained.

2.2 Restricted PCA

The second proposed method is named ‘‘Restricted PCA’’. A PCA on the original data matrix \mathbf{X} is computed to detect the Q relevant PCs (those with eigenvalue greater or equal to 1). Then each variable is assigned to the component with maximal weight, while all the other weights, for this variable, are set to zero. Finally, weights are column normalized. Formally:

- step 1 Compute the first Q PCs on Σ_X to obtain $\mathbf{Y} = \mathbf{X}\mathbf{A}^{(Q)}$; $\mathbf{B} = \mathbf{I}_Q$
- step 2 $\forall j$ compute: $m = \operatorname{argmax}\{|a_{jl}^{(Q)}|: l = 1, \dots, Q\}$, find component with largest loading, in absolute value, and retain it
- step 3 $a_{jq} = 0, q = 1, \dots, Q, q \neq m$, i.e., all the other weights are set to 0
- step 4 $\mathbf{A} = \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1/2}$ to normalize the weights.

2.3 Disjoint PCA

The third proposed method is named ‘‘DPCA’’. It is an extension of the model proposed by Vichi and Saporta (CDPCA, [9]) which focuses only on the classification of variables. In detail, DPCA classifies the variables into clusters (characterizing the corresponding components) by maximizing cluster-specific variance.

Given an $(n \times J)$ two-way two-mode (objects and variables) data matrix $\mathbf{X} = [x_{ij}]$, $i = 1, \dots, n$; $j = 1, \dots, J$ containing the measurements of J variables on n objects. Variables are supposed commensurate, if they are expressed in different scales they need to be standardized to have mean zero and unit variance.

The DPCA model can be formally written as follows

$$\mathbf{X} = \mathbf{X}\mathbf{A}\mathbf{A}' + \mathbf{E} \quad (1)$$

where \mathbf{A} is the component loading matrix with generally a reduced rank, i.e., $\text{rank}(\mathbf{A}) = Q \leq J$, satisfying constraints

$$\sum_{j=1}^J a_{jq}^2 = 1, \quad q = 1, \dots, Q \quad (2)$$

$$\sum_{j=1}^J (a_{jq}a_{jr})^2 = 0, \quad q = 1, \dots, Q-1; r = q+1, \dots, Q \quad (3)$$

and \mathbf{E} is an error component matrix. Note that $\mathbf{Y} = \mathbf{XA}$ specifies a reduced set of components and in this framework it can be assumed that these are corresponding to composite indicators of subsets of variables. Model (1) is the factorial model specifying the dimensionality reduction via the component loading matrix \mathbf{A} , which allows to partition variables into classes summarized by an orthonormal linear combination with maximal variance.

Moreover, the matrix \mathbf{A} is re-parametrized as the product of two matrices: \mathbf{B} and \mathbf{V} , where $\mathbf{V} = [v_{jq}]$ is a $(J \times Q)$ binary and row stochastic matrix defining a partition of variables into Q clusters, with $v_{jq} = 1$, if the j^{th} variable belongs to q^{th} cluster, $v_{jq} = 0$, otherwise; while, \mathbf{B} is a $(J \times J)$ diagonal matrix weighting variables such that

$$\sum_{j=1}^J v_{jq}b_j^2 = 1; \quad \sum_{q=1}^Q \sum_{j=1}^J v_{jq}b_j^2 = Q \quad (4)$$

Therefore, constraints (2) and (3) can be equivalently rewritten as

$$\mathbf{V} \text{ binary and row stochastic} \quad (5)$$

$$v_q' \mathbf{B}' \mathbf{B} v_q = 1, \quad q = 1, \dots, Q \quad (6)$$

Model (1) subject to constraints (2) and (3) can be considered the non-clustering version of the CDPCA.

3 SEM and PLS Path Modeling

Structural Equation Modeling is a methodology for representing, estimating, and testing a network of relationships between variables. As noticed in the introduction, we use the PLS-PM approach to estimate the SEM parameters. Specifically, it is composed by two sub-models: the measurement model (outer model) and the structural model (inner model).

The measurement model describes the relationships between the manifest and latent variables by using one of the following methods:

- the reflective method where each MV reflects its LV through a simple regression as described by

$$x_{jq} = \beta_{jq}\eta_q + \varepsilon_{jq}, \quad j = 1, \dots, J_q; \quad q = 1, \dots, Q \quad (7)$$

where

η_q is the q^{th} latent variable with $E(\eta_q) = 0$, $\sigma_{\eta_q} = 1$, $q = 1, \dots, Q$;

β_{jq} is the j^{th} weight on the q^{th} latent variable, $j = 1, \dots, J_q$; $q = 1, \dots, Q$;

ε_{jq} is the measurement error with $E(\varepsilon_{jq}) = 0$, $E(\varepsilon_{jq}, \eta_q) = 0$, $j = 1, \dots, J_q$; $q = 1, \dots, Q$;

- the formative method, where the LV is generated from its MVs through their linear combination, is given by

$$\eta_q = \sum_{j=1}^{J_q} w_{jq} x_{jq} + \delta_q \quad (8)$$

where

δ_q is the measurement error with $E(\delta_q) = 0$, $E(\delta_q, x_{jq}) = 0$, $j = 1, \dots, J_q$; $q = 1, \dots, Q$;

w_{jq} is the j^{th} weight on the q^{th} latent variable $j = 1, \dots, J_q$; $q = 1, \dots, Q$;

- the Mimic method which is a mixture of the previous two. The measurement model for a subset of MVs is the following

$$x_j = \beta_j\eta + \varepsilon_j; \quad j = 1, \dots, p \quad (9)$$

where the latent variable is defined by

$$\eta = \sum_{j=p+1}^{J_q} w_j x_j + \delta. \quad (10)$$

The first p MVs follow a reflective way and the $(J_q - p)$ remaining ones a formative way.

Note that, recently Esposito Vinzi and Russolillo [14] prefer to leave the measurement model always defined in the same way, i.e. as a factor model, no matter if formative or reflective scheme is used for outer estimation.

The latter sub-model is the structural model which describes the causal relationships between the L endogenous LVs as follows:

$$\eta_q = \sum_{l=1}^L \gamma_{sl}\eta_l + \zeta_q \quad (11)$$

where

γ_{sl} is the path coefficient that links the l^{th} to the s^{th} endogenous LV;

ζ_q is the residual with $E(\zeta_q) = 0$, $E(\zeta_q, \eta_q) = 0$, $q = 1, \dots, Q$.

Note that, the LVs considered into the SEM may be exogenous, that is, they have causes assumed to be external to the model, or endogenous, that is, they are predicted by the other variables in the model. To estimate the model parameters we use the PLS path algorithm. In detail, the initial step of the algorithm consists of randomly choosing an initial system of outer weights, then the LVs are estimated as both a linear combination of the MVs and a linear combination of its adjacent LVs. Finally, the outer weights are updated. The procedure is repeated until convergence is achieved, i.e. $\max \{w_{jq.ci} - w_{jq.pi}\} < \Delta$ where Δ is a convergence tolerance (usually $\Delta \leq 0.0001$).

4 Application to OECD Data Set

In this section, we apply the three different procedures, presented in Sect. 2, to the OECD data [12], in order to capture the potential underlying dimensions of well-being and to compare the three corresponding partitions. Moreover, we evaluate the relationships among those dimensions through a PLS-PM which allows to synthesize them into a unique final composite indicator. From a methodological point of view, the first issue for the construction of a composite indicator is the identification of variables which could adequately measure the potential underlying dimensions. The 20 variables used by OECD, plus 2 economic variables—observed on 34 member countries—are considered in this paper (see the first column in Table 1). They have been centered and scaled to unit variance in order to normalize different unit of measurements. By applying all the three proposed methods, we obtain 6 LVs (with variance greater than 1) which summarize different aspects of well-being. The structure of each LV is described in detail in Table 1. For all methods η_6 is defined through two economic variables (specifically x_{21} and x_{22}) and thus can represent a LV expressing *Economic well-being*. We can notice that the stepwise PCA and DPCA methods specify η_4 through the variables x_2 , x_{18} , and x_{20} which can represent a *Work & Life Balance* LV, and, moreover, approximately define also the other LVs in a similar way: η_5 can be defined as *Satisfaction* through x_1 , x_3 , x_5 , x_7 , x_{13} , x_{14} , and x_{15} for stepwise PCA, while through x_1 , x_3 , x_4 , x_7 , x_{13} , x_{14} , and x_{15} for DPCA; η_3 can be defined as *Social connections* through x_8 , x_9 , x_{10} , x_{12} , x_{16} , and x_{17} for stepwise PCA, while through x_8 , x_9 , x_{10} , x_{16} , and x_{17} for DPCA; η_2 can be defined as *Jobs* through x_4 , x_6 , and x_{19} for stepwise PCA, while through x_5 , x_6 , and x_{19} for DPCA; finally, η_1 can be defined as *Civic engagement* through x_{11} for stepwise PCA, while x_{11} and x_{12} for DPCA.

On the other hand, restricted PCA seems to have a slightly different behavior: x_1 , x_2 , x_3 , x_5 , x_7 , and x_{13} can represent η_5 as *Satisfaction and material assets*, x_6 , x_{14} , x_{15} , and x_{16} can define η_3 as *Social connections and personal care* while x_8 , x_{19} , x_{20} , and x_{11} define η_4 as representing *Work & Life Balance woman* LV, x_9 , x_{17} , x_4 , and x_{18}

Table 1 Components of the stepwise PCA (1), the restricted PCA (2), the disjoint PCA (3), and their PLS-PM communalities

	η_6			η_5			η_3			η_2			η_4			η_1		
	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3	1	2	3
x_{21} GDP pc	0.75	0.80	0.72															
x_{22} Consumer prices index	0.66	0.60	0.64															
x_1 Rooms per person				0.76	0.73	0.76												
x_3 Household disposable income				0.71	0.75	0.76												
x_5 Employment rate				0.65	0.66						0.94							
x_7 Quality of support network				0.70	0.68	0.67												
x_{13} Life expectancy				0.61	0.64	0.62												
x_{14} Self-reported health				0.48		0.52		0.58										
x_{15} Life satisfaction				0.73		0.72		0.73										
x_8 Educational attainment							0.44		0.43					0.55				
x_9 Students reading skills							0.68		0.69			0.71						
x_{10} Air pollution							0.46		0.47								0.71	
x_{16} Homicide rate							0.51	0.26	0.57									
x_{17} Assault rate							0.77		0.79			0.62						
x_4 Household financial wealth						0.39				0.46	0.25							
x_6 Long-term unemployment rate								0.35		0.45		0.39						
x_{19} Employment rate women child										0.60		0.80		0.55				
x_2 Dwelling without basic facilities					0.68								0.85		0.84			
x_{20} Time devoted pleasure personal													0.42	0.27	0.44			
x_{18} Employees working long hours											0.29		0.75	0.75				
x_{11} Voter turnout														0.06		1.00		0.11
x_{12} Consultation rule making							0.24										0.71	0.99

can represent η_2 as *Young Jobs* and, finally, x_{10} and x_{12} define η_1 as *Environmental Civic*.

Moreover, we postulate the same structural model for linking the LVs obtained by the three methods. Specifically it consists of four endogenous latent variables (η_6 , η_5 , η_4 , and η_3) and two exogenous latent variables (η_1 and η_2) as follows:

$$\eta_6 = \gamma_{16}\eta_1 + \gamma_{26}\eta_2 + \gamma_{36}\eta_3 + \gamma_{46}\eta_4 + \gamma_{56}\eta_5 + \zeta_6 \quad (12)$$

$$\eta_5 = \gamma_{15}\eta_2 + \gamma_{25}\eta_3 + \gamma_{35}\eta_4 + \zeta_5 \quad (13)$$

$$\eta_4 = \gamma_{14}\eta_2 + \gamma_{24}\eta_3 + \zeta_4 \quad (14)$$

$$\eta_3 = \gamma_{13}\eta_1 + \gamma_{23}\eta_2 + \zeta_3 \quad (15)$$

Note that the Eq. (12) specifies the effect of η_6 on all the remaining latent dimensions of well-being and it represents the composite indicator while the Eqs. (13), (14), and (15) represent the supposed linking between the LVs.

On the other hand, we specify a different measurement model (by using a reflective mode specification) according to the considered constrained PCA methods.

Specifically, for the stepwise PCA, it is as follows

$$x_{11,1} = \beta_{11,1}\eta_1 + \varepsilon_{11,1}$$

$$x_{4,2} = \beta_{4,2}\eta_2 + \varepsilon_{4,2}, x_{22} = \beta_{6,2}\eta_2 + \varepsilon_{6,2}, x_{19,2} = \beta_{19,2}\eta_2 + \varepsilon_{19,2}$$

$$x_{8,3} = \beta_{8,3}\eta_3 + \varepsilon_{8,3}, x_{9,3} = \beta_{9,3}\eta_3 + \varepsilon_{9,3}, x_{10,3} = \beta_{10,3}\eta_3 + \varepsilon_{10,3}, x_{12,3} \\ = \beta_{12,3}\eta_3 + \varepsilon_{12,3}, x_{16,3} = \beta_{16,3}\eta_3 + \varepsilon_{16,3}, x_{17,3} = \beta_{17,3}\eta_3 + \varepsilon_{17,3}$$

$$x_{2,4} = \beta_{2,4}\eta_4 + \varepsilon_{2,4}, x_{18,4} = \beta_{18,4}\eta_4 + \varepsilon_{18,4}, x_{20,4} = \beta_{20,4}\eta_4 + \varepsilon_{20,4}$$

$$x_{1,5} = \beta_{1,5}\eta_5 + \varepsilon_{1,5}, x_{3,5} = \beta_{3,5}\eta_5 + \varepsilon_{3,5}, x_{5,5} = \beta_{5,5}\eta_5 + \varepsilon_{5,5}, x_{7,5} = \beta_{7,5}\eta_5 \\ + \varepsilon_{7,5}, x_{13,5} = \beta_{13,5}\eta_5 + \varepsilon_{13,5}, x_{14,5} = \beta_{14,5}\eta_5 + \varepsilon_{14,5}, x_{15,5} = \beta_{15,5}\eta_5 + \varepsilon_{15,5}$$

$$x_{21,6} = \beta_{21,6}\eta_6 + \varepsilon_{21,6}, x_{22,6} = \beta_{22,6}\eta_6 + \varepsilon_{22,6};$$

for the restricted PCA, it consists in

$$x_{10,1} = \beta_{10,1}\eta_1 + \varepsilon_{10,1}, x_{12,1} = \beta_{12,1}\eta_1 + \varepsilon_{12,1}$$

$$x_{4,2} = \beta_{4,2}\eta_2 + \varepsilon_{4,2}, x_{9,2} = \beta_{9,2}\eta_2 + \varepsilon_{9,2}, x_{17,2} = \beta_{17,2}\eta_2 + \varepsilon_{17,2}, x_{18,2} \\ = \beta_{18,2}\eta_2 + \varepsilon_{18,2}$$

$$x_{6,3} = \beta_{6,3}\eta_3 + \varepsilon_{6,3}, x_{14,3} = \beta_{14,3}\eta_3 + \varepsilon_{14,3}, x_{15,3} = \beta_{15,3}\eta_3 + \varepsilon_{15,3}, x_{16,3} \\ = \beta_{16,3}\eta_3 + \varepsilon_{16,3}$$

$$x_{8,4} = \beta_{8,4}\eta_4 + \varepsilon_{8,4}, x_{11,4} = \beta_{11,4}\eta_4 + \varepsilon_{11,4}, x_{19,4} = \beta_{19,4}\eta_4 + \varepsilon_{19,4}, x_{20,4} \\ = \beta_{20,4}\eta_4 + \varepsilon_{20,4}$$

$$x_{1,5} = \beta_{1,5}\eta_5 + \varepsilon_{1,5}, x_{2,5} = \beta_{2,5}\eta_5 + \varepsilon_{2,5}, x_{3,5} = \beta_{3,5}\eta_5 + \varepsilon_{3,5}, x_{5,5} \\ = \beta_{5,5}\eta_5 + \varepsilon_{5,5}, x_{7,5} = \beta_{7,5}\eta_5 + \varepsilon_{7,5}, x_{13,5} = \beta_{13,5}\eta_5 + \varepsilon_{13,5}$$

$$x_{21,6} = \beta_{21,6}\eta_6 + \varepsilon_{21,6}, x_{22,6} = \beta_{22,6}\eta_6 + \varepsilon_{22,6};$$

and finally, for the DPCA, we have

$$x_{11,1} = \beta_{11,1}\eta_1 + \varepsilon_{11,1}, x_{12,1} = \beta_{12,1}\eta_1 + \varepsilon_{12,1}$$

$$x_{5,2} = \beta_{5,2}\eta_2 + \varepsilon_{5,2}, x_{6,2} = \beta_{6,2}\eta_2 + \varepsilon_{6,2}, x_{19,2} = \beta_{19,2}\eta_2 + \varepsilon_{19,2}$$

$$x_{8,3} = \beta_{8,3}\eta_3 + \varepsilon_{8,3}, x_{9,3} = \beta_{9,3}\eta_3 + \varepsilon_{9,3}, x_{10,3} = \beta_{10,3}\eta_3 + \varepsilon_{10,3}, x_{16,3} \\ = \beta_{16,3}\eta_3 + \varepsilon_{16,3}, x_{17,3} = \beta_{17,3}\eta_3 + \varepsilon_{17,3}$$

$$x_{2,4} = \beta_{2,4}\eta_4 + \varepsilon_{2,4}, x_{18,4} = \beta_{18,4}\eta_4 + \varepsilon_{18,4}, x_{20,4} = \beta_{20,4}\eta_4 + \varepsilon_{20,4}$$

$$x_{1,5} = \beta_{1,5}\eta_5 + \varepsilon_{1,5}, x_{3,5} = \beta_{3,5}\eta_5 + \varepsilon_{3,5}, x_{4,5} = \beta_{4,5}\eta_5 + \varepsilon_{4,5}, x_{7,5} \\ = \beta_{7,5}\eta_5 + \varepsilon_{7,5}, x_{13,5} = \beta_{13,5}\eta_5 + \varepsilon_{13,5}, x_{14,5} = \beta_{14,5}\eta_5 + \varepsilon_{14,5}, x_{15,5} \\ = \beta_{15,5}\eta_5 + \varepsilon_{15,5}$$

$$x_{21,6} = \beta_{21,6}\eta_6 + \varepsilon_{21,6}, x_{22,6} = \beta_{22,6}\eta_6 + \varepsilon_{22,6}.$$

In order to measure the quality of the measurement model, for each latent dimension, we compute the communalities, which measure the capacity of the path model to predict the MVs directly from their own LVs (see Table 1). As we can see, all values are quite high for all methods, meaning that the LVs are well identified by their MVs.

Figure 1 shows the path diagram of the estimated structural model obtained by PLS-PM with the stepwise PCA. As you can see, *Satisfaction* seems to be the most important driver for *Economic well-being* (path coefficient equal to 0.7508) followed by *Civic Engagement* and *Work and Life Balance*. Note that, *Social Connections* and *Jobs* seem to not have a significant effect on *Economic well-being* (p-value equal to 0.91 and 0.59, respectively). On the other hand, we can notice that significant relations exist also between *Jobs* and *Satisfaction* and between *Social connections* and *Work and life balance* (p-value equal to 0.00001 and 0.0118, respectively).

Similar results have been obtained by PLS-PM with restricted PCA, as shown in Fig. 2. In fact, also in this case, *Satisfaction* has the highest effect on *Economic well-being* with a path coefficient equal to 0.846. The main differences stay on the effect of *Social connections* on *Satisfaction* and *Jobs* on *Work and life balance*, which are greater with respect to the previous analysis.

Finally, the path diagram for the structural model obtained by using the third procedure (DPCA) is shown in Fig. 3. As you can observe, similar interpretation

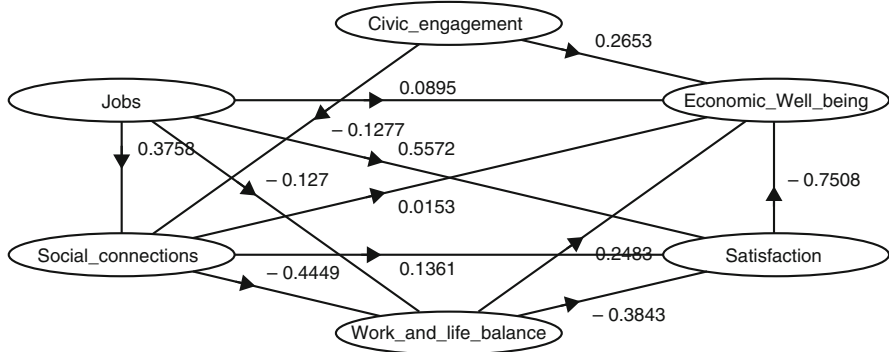


Fig. 1 Path diagram of the structural model obtained by using PLS-PM with stepwise PCA

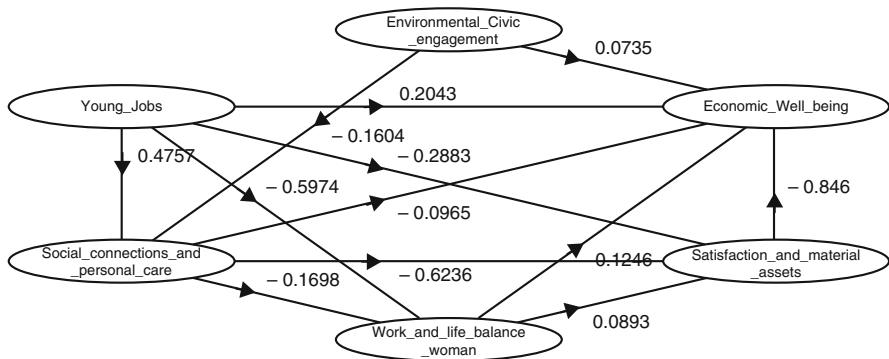


Fig. 2 Path diagram of the structural model obtained by using PLS-PM with restricted PCA

with respect to the previous procedures can be done. In fact, we have that *Satisfaction* has the highest effect on *Economic well-being* followed by *Work and Life Balance*.

Moreover, we show in Table 2, the estimated path coefficients related to the three composite indicators.

As you can observe, while for the stepwise PCA η_1 has a significant effect on η_6 , for the restricted PCA and the DPCA this effect is negligible. On the other hand, for the restricted PCA η_2 seems to show an effect on η_6 which is absent in the other two procedures. For all methods, while η_3 seems to not have any effect on η_6 , η_4 shows a discrete impact on η_6 , and η_5 is the most important driver.

In Fig. 4, we show countries ranking according to the composite indicators obtained by means of each of the three proposed methods. In particular, we can observe that the lower positions are generally stable for all methods, while for the higher position we have more variability among the three analyses.

Finally, to compare proposals, we use two benchmark measures: the GoF index (an overall assessment of the model) and the proportion of explained variance of the MVs (Table 3).

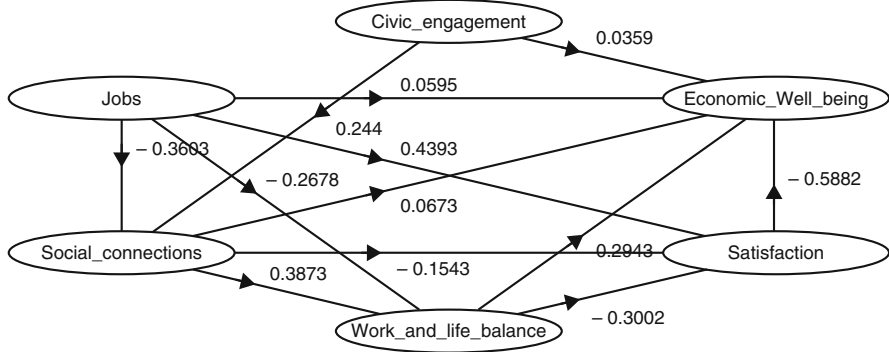


Fig. 3 Path diagram of the structural model obtained by using PLS-PM with DPCA

Table 2 Estimated path coefficients

	η_6 stepwise PCA	η_6 restricted PCA	η_6 DPCA
η_1	0.2653	-0.0735	0.0358
η_2	-0.0895	0.2043	-0.0594
η_3	-0.0153	-0.0965	0.0673
η_4	0.2483	-0.1246	0.2943
η_5	0.7508	0.8460	0.5882

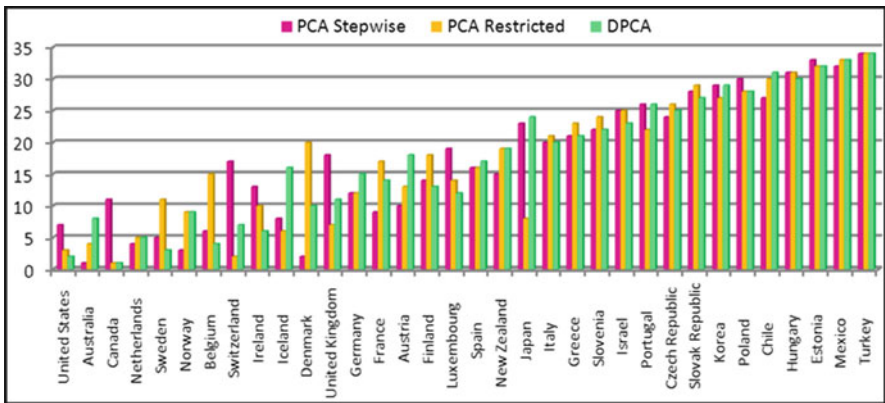


Fig. 4 Plot of the rankings obtained by using three different approaches

Table 3 Benchmark measures

Methods	Proportion of explained variance	GoF PLS-PM
Stepwise PCA	0.832	0.749
Restricted PCA	0.653	0.818
Disjoint PCA	0.863	0.755

Thus, if a composite indicator has to be defined it has to explain the largest part of the variance of the original variables. Hence, the best model is specified by the DPCA. However, if prediction of well-being is an important issue, restricted PCA defines the most suitable model.

5 Conclusions and Challenges for Further Research

In this paper, we have proposed three different constrained version of PCA to capture the latent underlying dimensions of well-being. Furthermore, after forming the latent variables, we study their relationships by using the SEM under PLS-PM approach. These methodologies are compared by using a data set from 34 member countries of the OECD. Notice that, the use of SEMs is particularly appropriate in these kinds of studies where simultaneously the aim is to analyze the relationships between several variables (economic and not-economic) and potential latent dimensions related to well-being. In particular, in this work, we use a sequential approach, consisting of building the latent dimensions by using an aggregation method, and then we give its results in input to the confirmatory SEM. In future researches, we would like to propose a unique model able to simultaneously identify latent dimensions, and estimate their connections. It would be a good research line also to extend the PLS-PM to longitudinal studies.

References

1. Saisana, M., Tarantola, S.: State-of-the-art Report on Current Methodologies and Practices for Composite Indicator Development. EUR 20408 EN, European Commission-JRC (2002)
2. Kline, R.B.: Principles and Practice of Structural Equation Modeling. The Guilford Press, New York (2011)
3. Jöreskog, K.: A general approach to confirmatory maximum likelihood factor analysis. *Psychometrika* **34**, 183–202 (1969)
4. Thurstone, L.L.: Multiple-Factor Analysis. University of Chicago Press, Chicago (1947)
5. Reiersol, O.: On the identifiability of parameters in Thurstone's multiple factor analysis. *Psychometrika* **15**, 121–149 (1950)
6. Zou, H., Hastie, T., Tibshirani, R.: Sparse principal component analysis. *J. Comput. Graph. Stat.* **15**(2), 265–286 (2006)
7. Tenenhaus, M., Esposito Vinzi, V.: PLS regression, PLS path modeling and generalized procrustean analysis: a combined approach for PLS regression, PLS path modeling and generalized multiblock analysis. *J. Chemom.* **19**, 145–153 (2005)
8. Tenenhaus, A., Tenenhaus, M.: Regularized generalized canonical correlation analysis. *Psychometrika* **76**(2), 257–284 (2011)
9. Vichi, M., Saporta, G.: Clustering and Disjoint Principal Component. *Comput. Stat. Data Anal.* **53**(8), 3194–3208 (2009)
10. Bollen, K.A.: Structural Equations with Latent Variables. Wiley, New York (1989)
11. Kaplan, D.: Structural Equation Modeling: Foundations and Extensions. Thousands Oaks, Sage (2000)
12. OECD: How's Life?: Measuring Well-Being. OECD Publishing (2011)

-
13. Hall, J., Giovannini, E., Morrone, A., Ranuzzi, G.: A Framework to Measure the Progress of Societies. OECD Statistics Working Papers, No. 2010/05, OECD, Paris (2010)
 14. Esposito, V.V., Russolillo, G.: Partial least squares algorithms and methods. Wiley Interdiscip. Rev. Comput. Stat. **5**, 1–9 (2013)

Extracting Meta-information by Using Network Analysis Tools

Agnieszka Stawinoga, Maria Spano, and Nicole Triunfo

Abstract

This paper has been developed in the frame of the European project BLUE-ETS (Economic and Trade Statistics), in the work-package devoted to propose new tools for collecting and analyzing data. In order to obtain business information by documentary repositories, we refer to documents produced with nonstatistical aims. The use of secondary sources, typical of data and text mining, is an opportunity not sufficiently explored by National Statistical Institutes. The use of textual data is still viewed as too problematic, because of the complexity and the expensiveness of the pre-processing procedures and often for the lack of suitable analytical tools. In this paper we pay attention to the problems related to the pre-processing procedures, mainly concerning with semantic tagging. We propose a semi-automatic strategy based on network analysis tools to create financial-economic meta-information useful for the semantic annotation of the terms.

1 Introduction

This work has been developed in the frame of the European project BLUE-ETS, acronym for BLUE Enterprise and Trade Statistics (www.blue-ets.istat.it), funded by the European Commission (7th Framework Programme). Our peculiar task in BLUE-ETS consists in proposing new tools for collecting and analyzing data. In order to reduce response burden and, at the same time, to collect cheaper and (better) quality data, we refer to secondary sources, produced with nonstatistical aims. The use of secondary sources, typical of data and text mining, is an opportunity not

A. Stawinoga • M. Spano (✉) • N. Triunfo
University of Naples Federico II, Napoli, Italy
e-mail: agnieszka.stawinoga@unina.it; maria.spano@unina.it; nicole.triunfo@unina.it

sufficiently explored by National Statistical Institutes. NSIs aim at collecting and representing information in a usable and easy-readable way. The use of textual data has been still viewed as too problematic, because of the complexity and the expensiveness of the pre-processing procedures and often for the lack of suitable analytical tools. In order to focus our attention on business statistics we propose to extract and to analyze data by mining into the management commentaries attached to the annual reports of the companies listed on Italian Stock Exchange. In this work we pay attention to the problems related to the pre-processing procedures, mainly concerning with semantic tagging. In order to perform the semantic tagging it is necessary to have a language resource appropriate for the subject of the analysis. In this paper we propose a semi-automatic strategy based on network analysis tools for creating financial-economic meta-information.

2 Pre-processing: State of Art

Statistical studies on natural language have been developed thanks to the evolution of computer resources so as to produce the automatic analysis of texts [10]. The increasing availability of computerized linguistic resources [15] and the increasing popularity of texts and documents available on-line have modified the criteria and the techniques giving the possibility of a further use of this kind of data. The solutions are based on a multidisciplinary approach combining statistical and linguistic tools in the text mining context [7, 13, 16]. Text mining aims at identifying in an automatic or semi-automatic way the regularities existing in large textual databases. Textual data analysis methods give the possibility to extract useful and not trivial information from large databases. The main problem of automatic text analysis is to understand and to extract the real meaning of the documents. In order to deal with this issue, it is necessary to transform textual (unstructured) data in a lexical matrix which can be analyzed with statistical tools. This transformation concerns the selection of the terms which are able to represent the syntagmatic axis of the document. In the literature this process is well-known as text pre-processing. A unique definition of the pre-processing steps does not exist. According with the aim of the analysis, the researcher creates an ad hoc strategy for extracting the significant information from the documents. The researchers' choices affect the results of the analysis. In order to perform an automatic text analysis, the first step consists in choosing the unit of the analysis. On the one hand, the formalists [4, 11, 12] consider the graphical form/type (sequence of characters delimited by two separators) as the unit of analysis. They carry out the statistical analysis regardless of the meaning of the graphical forms, which are considered language independent. On the other hand, the computational linguists consider the lemma as the unit of the analysis [8]. Electronic dictionaries, frequency lexica, and automatic normalization are the practical tools adopted according with this language dependent approach. In the field of textual data analysis, Bolasco [5] considers a mixed language dependent unit (graphical form/lemma/multiword expression)

named “textual form.” Once the unit of analysis has been chosen, the pre-processing can be summarized in the following steps:

1. cleaning of the text;
2. normalization;
3. text annotation.

The cleaning of the text consists of the definition of alphabet characters/separators. The normalization consists of the recognition of particular entities such as dates, numbers, currencies, acronyms, abbreviations, place names, and phrases or expressions of interests. Another part of this preliminary phase is the text annotation. It consists in associating meta-information (grammatical and semantic tagging, lemmatization, etc.) with terms [6]. The semantic tagging is very useful for identifying important terms in subsequent phases of the analysis. It enriches the unstructured or semi-structured data with a context which is further linked to the structured knowledge of a domain, such as the economic and financial.

3 Methodology

In this paper we propose to extend the strategy for detecting meaningful communities of closely related terms introduced by Balbi and Stawinoga [2]. One of the principal ideas of the strategy is to represent the textual data in the form of a network. The data matrix used for the analysis is the lexical table \mathbf{T} ($n \times p$) which indicates the occurrences of the p terms in n parts (documents) of a corpus. For the aim of the analysis, the documents are considered as the units and the terms as the variables. The matrix \mathbf{T} can be easily transformed into a one-mode co-occurrence matrix \mathbf{W} by the product $\mathbf{A}^t\mathbf{A}$. \mathbf{A} is a binary matrix where the generic element $a(i,j)$ ($i = 1, \dots, n; j = 1, \dots, p$) equals 1 if the term j occurred at least once in the document i and 0 otherwise. The adjacency valued matrix \mathbf{W} represents the relational system of the selected terms. To avoid weak relations and to obtain less sparse matrix the authors proposed to use the strength of associations among terms instead of the simply co-occurrence frequency. To extract the graphs of similarities they applied the Jaccard index. The similarity between two terms k and j is defined as:

$$S_{kj} = \frac{w_{kj}}{w_{kk} + w_{jj} - w_{kj}}. \quad (1)$$

According to a predefined threshold based on the actual distributions of Jaccard index, the matrix of similarities \mathbf{S} ($p \times p$) is dichotomized to obtain network representation of the relations existing among the terms. This strategy has been applied to the management commentary of the world leader of eyewears, Luxottica, because this company is listed on both the US and Italian markets and the US law

thoroughly explains all information that the management commentary must include (differently to the Italian law).

In this paper we propose to extend this strategy to the analysis of more than one management commentary, in order to identify meta-information (statistical language resources) for the financial-economic field. In this case, the rows of the lexical table **T** represent the different companies and the columns are the terms. Once the relational structure of the terms is detected by using different network tools we go to identify topics which are mostly treated by all companies in their annual reports. The network analysis tools, we propose to use, are the following. Firstly, different components of the network (maximal connected sub-graphs) are extracted. From a textual point of view, the different connected components give the possibility to individuate different topics.

When a component of the network is too big to identify a single concept (e.g., the main component), we propose to analyze the cohesion of this sub-graph. Firstly, we go to calculate the density which is simply the proportion of ties actually present and the total number of possible ties in the network. This measure varies from 0 (empty graph) to 1 (complete graph) and its high value indicates high cohesion of the graph. On the other hand, an important macro-characteristic of the network are the distances among nodes. The geodesic distance between two nodes is the length of any shortest path between them [14]. The average geodesic distance between all actors in a graph indicates how close actors are together. The smaller its value the higher cohesion of the network is observed. The longest geodesic distance (diameter) can be considered as an alternative measure to control the level of cohesion of the (connected) network and it quantifies how far apart the farthest two nodes in the graph are. The index distance-based cohesion (compactness) measures the normalized sum of the reciprocal of all the distances. This index ranges from 0 (the network is entirely made up of isolates) to 1 (everyone is adjacent in the network) so large values indicate high cohesiveness of the network.

For the main component with low level of cohesion we propose to calculate betweenness centrality [9] to individuate the most influential nodes connecting different communities within this sub-graph of the network. The betweenness centrality index measures how often a node appears on the shortest path between two other random nodes in the network. This measure is defined as:

$$C_B(i) = \sum_{j < k} \frac{g_{jk}(i)}{g_{jk}}, i \neq j \neq k \quad (2)$$

where $g_{jk}(i)$ is the number of shortest paths between nodes j and k that contain the node i , g_{jk} represents the number of shortest paths between nodes j and k . In terms of text analysis, the higher the value, the more often the term links different contexts within the text. In the network analysis of dictionaries, Batagelj et al. [3] use betweenness centrality index to individuate terms which are used as intermediates for explaining other terms. As for large-scale networks calculation of the betweenness centrality is often very heavy from a computational viewpoint,

a parallel algorithm proposed by Bader and Madduri [1] gives the possibility to deal with this problem. The proposed strategy allows us to pass from elementary data (terms) to higher order data (context, topics). It gives the possibility to select a subset of relevant terms, which illustrate important topics characterizing the corpus.

4 Results

We analyze a sample of 50 Italian listed companies, extracted using a technique of quota sampling to ensure compliance with the composition of the sectors in which the firms in the Italian Stock Exchange (Borsa Italiana s.p.a.) are classified. We exclude financial firms because their management commentary is subjected to a specific regulation and it could be very different from those of non-financial firms. Our reference year is 2010.

Only a part of the management commentaries, commonly known as “results of operations,” has contributed to build the textual database. The pre-processing procedures are performed using the software TalTac 2.0. This step gives the possibility to clean up the text from empty words (conjunction, articles, adverbs, etc.) and words with frequency < 5 (rare words). Following this procedure, we managed to get the lexical matrix of size (50×541) . Applying the different steps of the strategy described in Sect. 3 we obtain the similarity matrix \mathbf{S} (541×541) . According to the actual distribution of Jaccard index in the data, we choose a threshold value 0.55 to dichotomize the matrix \mathbf{S} . In this way we obtain the adjacency matrix \mathbf{X} which represents the network of the terms (541 nodes and 774 edges).

This network consists of 371 isolates (nodes which are not connected to any other nodes), 32 components of 2 nodes, 7 components of 3 nodes, 2 components of 5 nodes, 2 components of 9 nodes, and 1 component of 61 nodes. In order to illustrate important topics characterizing the corpus we consider only components of size greater than or equal to 2 (see Fig. 1).

The main (the biggest) component (see Fig. 2) of the network consists of 61 nodes and 694 edges.

From a statistical viewpoint it is a sub-graph strongly connected as confirmed by the main indicators of compactness: the distance-based cohesion, the average distance, the density, and the diameter. The distance-based cohesion (compactness) of the component is, in this case, equal to 0.666. This index varies between 0 and 1, so high values indicate greater cohesion. Furthermore, the average geodesic distance between all pairs of nodes is represented by a low value, equal to 1.769. This value is inversely proportional to the density of the sub-graph which is equal to 0.37. In the current case, the value of diameter indicates that no actor is more than four steps from any other so we can consider the main component as a compact sub-graph. The high compactness of the main component is confirmed also by the values of betweenness centrality index which indicate the lack of the most influential nodes connecting different communities within this sub-graph of the network. From a textual viewpoint, the high cohesion of the main component is due to the presence

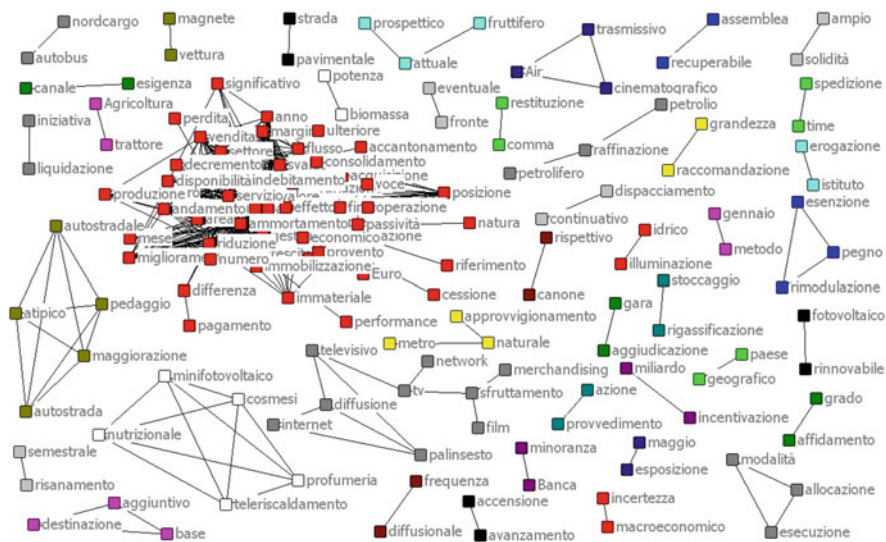


Fig. 1 Components extracted from the analyzed network

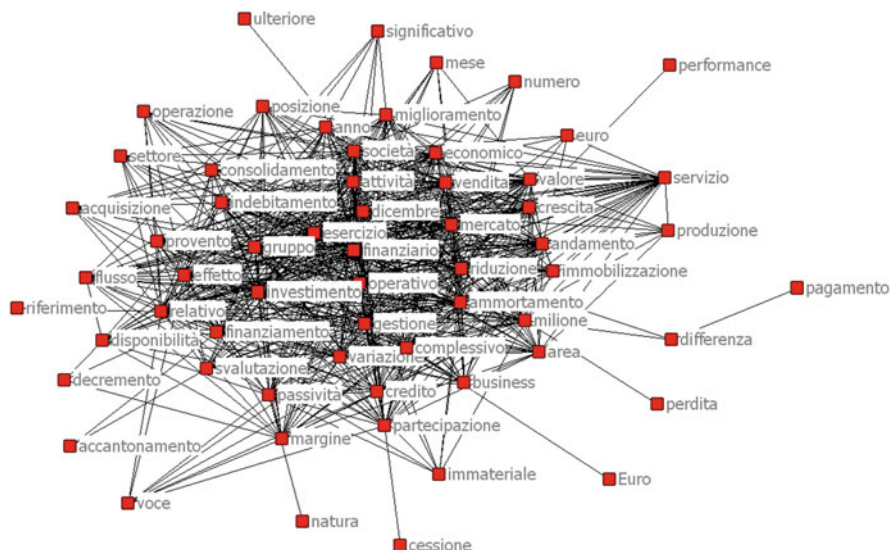


Fig. 2 The main component of the analyzed network

of terms used by all companies in order to illustrate their results of operations. The drafting of the financial statements involves the use of a specific vocabulary, such as the economic and financial lexicon. The terms of the main component can be used to build a statistical language resource in order to semantically tag the terms within the corpus. According to the aim of the researcher this group of terms, which are

common for all documents, could be also a stop list useful for deleting the “trivial” terms.

Compared to the main component, the components of 2, 3, 5, and 9 nodes represent groups of terms which characterize only some of the analyzed documents. Therefore they are separated from the rest of the network. Specifically, each component expresses a conceptual atom which gives the possibility to identify homogeneous groups of documents (companies), such as it is shown in Fig. 3. For instance, the component *raffinazione–petrolio–petrolifero* represents, respectively, the process, the product, and the market of the companies (Saras, Eni, Erg) operating in the energy sector. One of the companies (Erg) is connected through an arc to another component *rinnovabile–fotovoltaico*, consisting of two nodes, which represents the other business area “green economy” in which this company operates. Compared to the companies linked to the latter component, Kinexia stands out because it is characterized by a group of terms *minifotovoltaico–profumeria–cosmesi–teleriscaldamento–nutrizionale*. This component represents the non-characteristic activity of Kinexia. One of the components of five nodes *autostrada–atipico–maggiorazione–pedaggio–autostradale* expresses the event which has characterized the operating results of the companies (Atlantia, Autostrade Meridionali) operating in the field of highway infrastructure. In 2010 this sector recorded an increase in net operating incomes generated by an increase in tolls.

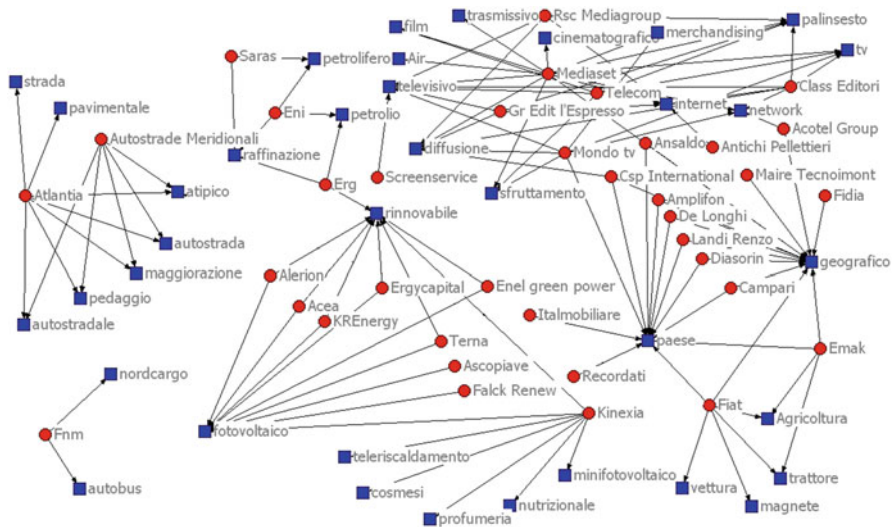


Fig. 3 The two-mode network of the extracted components

5 Conclusions and Future Work

In this paper we propose a semi-automatic strategy to extract useful information starting from administrative documents expressed in natural language. The strategy can be summarized in the following steps:

1. cleaning of the textual database;
2. representation of the textual data in the form of a network;
3. exploration of the relational structure existing among the terms by using network analysis indexes;
4. identification of the network components.

The proposed strategy makes it possible to identify the components (groups of terms), which characterize the analyzed documents. These components can be used:

1. to select the most relevant terms viewed as new features prior to further analyses;
2. to identify the different topics treated by the companies in the drafting of their administrative documents;
3. to reduce the dimensionality by passing from elementary data (terms) to higher order data (context, topics).

Further developments of the research will be devoted to put the strategy into the theoretical framework of dimensionality reduction. The research will be centered on tools for synthesizing elementary data in higher order data. Another issue worth to be investigated is the choice and comparison of metrics for building the similarity matrix in order to fully explore relational structures among the terms in a corpus.

Acknowledgements This work is financially supported by the European Project BLUE-ETS.

This paper derives by a strict and continuous collaboration among the authors. Anyway Sects. 1 and 4 may be mainly attributed to M. Spano; Sect. 3 to A. Stawinoga; Sects. 2 and 5 to N. Triunfo.

References

1. Bader, D.A., Madduri, K.: Parallel algorithms for evaluating centrality indices in realworld networks. In: Proceedings of the 35th International Conference on Parallel Processing (ICPP). IEEE Computer Society, Columbus, OH (2006)
2. Balbi, S., Stawinoga, A.: The use of Network Analysis tools for dimensionality reduction in Text Mining, SLDS 2012, Florence. <https://www.docenti.unina.it/ricerca/visua/lizzaAttivitaRicerca.do?idDocente=53494d4f4e4142414c4249424c42534d4e35384c35394638333944&nomeDocente=SIMONA&cognomeDocente=BALBI> (2012)
3. Batagelj, V., Mrvar, A., Zaveršnik, M.: Network analysis of dictionaries. In: Erjavec, T., Gros, J. (eds.) *Jezikovne tehnologije /Language Technologies*, Ljubljana, pp. 135–142 (2002)
4. Benzécri, J.P.: *Pratique de l'Analyse Des Données, Linguistique e Lexicologie*. Dunod, Paris (1981)

5. Bolasco, S.: Sur différentes stratégies dans une analyse des forms textuelles: une expérimentation à partir de données d'enquête. In: Bécue, M., Lebart, L., Rajadell, N. (eds.) JADT 1990, UPC, Barcelona, pp. 69–88 (1990)
6. Bolasco, S.: Meta-data and strategies of textual data analysis: problems and instruments. In: Hayashi et al. (eds.) Data Science, Classification and Related Methods (Proceedings V IFCS - Kobe, 1996). Springer, Tokio, pp. 468–479 (1998)
7. Bolasco, S., Canzonetti, A., Capo, F.: Text Mining: Uno strumento strategico per imprese e istituzioni. Cisu Editore, Roma (2005)
8. De Mauro, T.: I vocabolari ieri e oggi. In: Il vocabolario del 2000, a cura di IBM Italia, Roma (1989)
9. Freeman, L.C.: Centrality in social networks conceptual clarification. *Soc. Netw.* **1**, 215–239 (1979)
10. Lebart, L., Salem, A.: Statistique textuelle. Dunod, Paris (1994)
11. Reinert, M.: Un logiciel d'analyse lexicale: ALCESTE. *Les Cahiers de l'analyse des données*, **XI 4**, 471–484 (1986)
12. Salem, A.: Pratique des segments répétés. Essai de statistique textuelle. Klincksieck, Paris (1987)
13. Sullivan, D.: Document Warehousing and Text Mining: Techniques for Improving Business Operations, Marketing and Sales. Wiley, New York (2001)
14. Wasserman, S., Faust, K.: Social Network Analysis: Methods and Applications. University Press, Cambridge (1994)
15. Zampolli, A., Calzolari, N.: Problemi, metodi e prospettive nel trattamento del linguaggio naturale: l'evoluzione del concetto di risorse linguistiche. In: Cipriani, R., Bolasco, S. (eds.), pp. 51–68 (1995)
16. Zanasi, A.: Text Mining and Its Applications to Intelligence, CRM and Knowledge Management. WIT Press, Southampton (2005)

Factor PD-Co-clustering on Textual Data

Cristina Tortora, Marina Marino, and Germana Scepi

Abstract

In this paper we propose to extend factor probabilistic distance (FPD) clustering to FPDco-clustering for frequency data. FPD-clustering transforms the data using a factor decomposition and clusters the transformed data optimizing the same criterion. FPDco-clustering simultaneously finds clusters of rows and column basing on the PD-clustering criterion. The method is useful in case of large data sets. In this paper the new method is applied on large textual data sets with the aim of extracting interesting information.

1 Introduction

Clustering aims to group subjects with similar characteristics. Standard clustering techniques, also called one-way clustering techniques, are based on similarities among subjects across all variables. In some fields, such as genetics, biology or

Proceeding of the 46th Scientific Meeting of the Italian Statistical Society.

C. Tortora
McMaster University, Hamilton, ON, Canada
e-mail: ctortora@math.mcmaster.ca

M. Marino
Dipartimento di Scienze Sociali, Università degli Studi di Napoli Federico II, Napoli, Italy
e-mail: mari@unina.it

G. Scepi (✉)
Dipartimento di Scienze Economiche e Statistiche, Università degli Studi di Napoli Federico II,
Napoli, Italy
e-mail: scepi@unina.it

textual analysis, clustering both units and variables may be of interest. The original idea is to find clusters of similar elements showing similar clusters of variables. The first reference to both clustering of rows and columns is [11] proposing what was later called a “sequential approach”. That approach clusters successively and independently the rows and the columns of the data matrix. Otherwise, simultaneous clustering algorithms seek “blocks” of rows and columns that are interrelated. There are many advantages in a simultaneous rather than one-way clustering or a sequential approach [2]. Simultaneous clustering highlights the association between the clusters of rows and columns, it allows the researcher to deal with sparse and high-dimensional data matrices and has peculiar advantages when we deal with unsupervised data analysis. Van Mechelen et al. [23] and Charrad and Ahmed [3] show a review of simultaneous clustering techniques. In the literature many approaches are based on a simultaneous rather than a sequential approach [2,6,9,10,13,15,17,23,24]. This type of clustering approach has been defined with a big range of different names like: direct clustering, block clustering, bi-dimensional clustering, two-way clustering, co-clustering or bi-clustering, each name refers to specific characteristic of the method. When the cluster of one dimension depends on the cluster of the other dimension the method is defined as two-way clustering [7, 8, 16, 19], among them [4] uses the term co-clustering when referring to textual data two-way clustering. Because of our focus on textual data set and because of the characteristic of the method we propose, hereafter, we will use the term co-clustering. The aim of this paper is to propose a new factor co-clustering method based on factor probabilistic distance clustering (FPDC) [22]: FPDco-clustering (FPDcoC). FPDC, and consequently FPDcoC, has the advantage that can be used with high-dimensional and sparse matrices, even with more rows than columns and it can deal with non-spherical clusters, outliers or noisy data. Co-clustering methods are often used on textual data sets [1, 5]. Texts can be transformed in a words by texts data matrix X ; each cell x_{ij} indicates the frequency of the i th word in the j th text. Texts may be clustered based upon their word distributions or words may be clustered based upon their distribution among texts. However, it is often desirable to co-cluster or simultaneously cluster both dimensions of a contingency table by exploiting the clear duality between rows and columns. Moreover, as Dhillon et al. [5] suggests, even if we are interested in clustering along one dimension of the contingency table, when dealing with sparse and high-dimensional data (as it is always the case when we deal with word-document matrices), it turns out to be beneficial to employ co-clustering approaches. The paper has the following structure. In Sect. 2 we present factor PD-clustering method and its extension for frequency data sets. The new method FPDcoC is shown in Sects. 3 and 4 shows a simulation study and in Sect. 5 FPDcoC is applied on a real textual data set. Section 6 shows some conclusions.

2 Background: Factor PD-Clustering

Let us define with X , an $n \times J$ data matrix, where x_i is a J -dimensional vector. Given K clusters that are assumed to not be empty, let's define with $d_k(x_i)$ the distance of a point x_i from the cluster the J -dimensional centre c_k and with p_{ik} the probability of a point x_i to belong to the cluster k , with $i = 1, \dots, n$ and $k = 1, \dots, K$. PD-clustering basic assumption is that the product between the probability of any point belonging to a cluster and the distance from the centre of the cluster is a constant, $p_{ik}d_k(x_i) = F(x_i)$, $\forall i = 1, \dots, n$ and $\forall k = 1, \dots, K$.

The value of the constant $F(x_i)$ is a measure of the closeness of x_i to the cluster centres, it measures the classifiability of x_i , where a point x_i is more classifiable than a point x_j with respect to a centre c_k if $d_k(x_i) < d_k(x_j)$. $F(x_i)$ is a maximum when the distance between x_i and all the cluster centres is the same, it is close to zero if x_i is close to one of the cluster centre. The sum over i of $F(x_i)$ is called join distance function, JDF . The aim of the algorithm is to maximize the classifiability of all points. The whole clustering problem consists in the identification of the centres that minimize the JDF. Basing on this assumption, cluster centres can be computed as follows:

$$c_k = \sum_{i=1, \dots, N} \left(\frac{u_k(x_i)}{\sum_{j=1, \dots, N} u_k(x_j)} \right) x_i,$$

where $u_k(x_i) = \frac{p_{ik}^2}{d_k(x_i)}$ (for details refer to [12]).

When the number of variables is large and variables are correlated non-hierarchical clustering methods, including PD-clustering, become very unstable and the correlation between variables can hide the real clustering structure [25]. A linear transformation of the original variables into a reduced number of orthogonal ones can significantly improve the algorithm performance. The linear transformation of the variables and the clustering method must optimize a consistent criterion. Let define with G a distance array of general element $g_{ijk} = |x_{ij} - c_{kj}|$, where $i = 1, \dots, n$ indicates the units, $j = 1, \dots, J$ the variables and $k = 1, \dots, K$ the clusters. In [20] it is demonstrated that a Tucker3 decomposition [14] of the distance array G minimizes the JDF, or equivalently optimizes the same criterion of PD-clustering. The method consists in finding a transformation of the original data $x_{iq}^* = x_{ij}b_{jq}$, where b_{jq} is a weighting system, and cluster centres c_{kq}^* such that the JDF is minimized, where x_{iq}^* and c_{kj}^* indicate the projection of points and centres in the factor space, C and B are the matrix of centres and the weights, respectively. An iterative algorithm is used to obtain x_{iq}^* and c_{kj}^* .

1. random initialization of matrix C and computation of G ;
2. Tucker3 decomposition of G to obtain x_{iq}^* ;
3. computation of c_{kj}^* using PD-clustering and update G .
4. if JDF decreases, go to step 2, stop elsewhere.

The convergence of the method is empirically demonstrated. The integration of the PD-clustering and the Tucker3 step makes the clustering more stable and allows one to consider data sets with large number of variables and not elliptically shaped [22].

When the data represent frequencies the previously defined distance cannot be used. The most used distance for frequency matrix is the χ^2 . The χ^2 distance among two generic points x_i and x_m is defined as in the following expression:

$$d(x_i, x_m)_{\chi^2} = \frac{\sum_{j=1}^J \frac{(x_{ij} - x_{mj})^2}{(x_{ij} + x_{mj})}}{2}. \quad (1)$$

Using the (1) the JDF becomes:

$$\text{JDF} = \sum_{i=1}^n \sum_{k=1}^K d(x_i, c_k)_{\chi^2} p_{ik}^2 \quad \text{s.t.} \quad \sum_{i=1}^n \sum_{k=1}^K p_{ik}^2 \leq n, \quad (2)$$

and the objective function is

$$\text{JDF}^* = \min_{C,B} \sum_{i=1}^n \sum_{k=1}^K d(x_{ij} b_{jq} - c_{kj} b_{jq})_{\chi^2} p_{ik}. \quad (3)$$

In that case the weighting system that minimizes the JDF can be found through a singular value decomposition of the matrix D of elements $d(x_i, c_k)_{\chi^2}$.

3 Factor PD-Co-clustering

FPDC can be extended to a two-mode clustering, FPDcoC. A simple application of FPDC on the row and the column of the matrix will lead to two independent clustering, one for the row and one for the column. The result would be a centre matrix C and a weighting system B for the units and centre matrix $C^{(v)}$ of generic element $c_k^{(v)}$ and a weighting system $B^{(v)}$ for the variables. The aim of two-mode clustering is to cluster the data matrix simultaneously [23]. Let define with D a distance matrix of order $n \times K$ with general element $d(x_i, c_k)_{\chi^2}$, with B the singular vectors of D , with y_j the general element of the transposed data matrix $X^{(t)}$, with $D^{(v)}$ a distance matrix of order $P \times K$ with general element $d(y_j, c_k^{(v)})_{\chi^2}$, with $B^{(v)}$ the singular vectors of $D^{(v)}$. A singular value decomposition (svd) of D allows to project the variables y^* on the factor space that maximizes distances among units: $y_{ji}^* = y_{ji} b_{ij}$, with $i = 1, \dots, n, j = 1, \dots, J$. On the projection obtained, variables y^* can be grouped using PD-clustering. Specifically, one-mode clustering optimizes the (3), while two-mode clustering optimizes

$$\text{JDF}^* = \min_{C^{(v)}, B^{(v)}} \sum_{i=1}^n \sum_{k=1}^K d(y_{ji}^* b_{ii}^{(v)} - c_{ki}^{(v)} b_{ii}^{(v)}) \chi^2 p_{ik}, \quad (4)$$

where $b_{ii}^{(v)}$ is the general element of $B^{(v)}$. Starting from the variable clustering structure a distance matrix of variables $D^{(v)}$ can be computed. Basing on the same procedure, units are projected on the space that maximizes the classifiability of variables, $x_{ij}^* = x_{ij} b_{jj}^v$, the clustering partition is obtained in the new space. The objective function is

$$\text{JDF}^* = \min_{C, B} \sum_{i=1}^n \sum_{k=1}^K d(x_{ij}^* b_{jj}, c_{kj} b_{jj}) \chi^2 p_{ik}^2. \quad (5)$$

The entire process is iterated until the convergence is reached.

FPDcoC can be summarized in the following steps:

1. Random initialization of cluster of units and computation of distance matrix D ;
2. svd of distance matrix D and projection of variables y_{ji}^* ;
3. PD-clustering of variables on the factor space;
4. Computation and svd of distance matrix of variables D^v , projection of units x_{ij}^* ;
5. PD-clustering of units on the factor space;
6. if JDF decreases, go to step (2), stop elsewhere.

The convergence is empirically demonstrated. An advantage of the method is that it is not based on the arithmetic mean and it is not affected by the presence of outliers in the data [21], however the method catches small clusters if there are. The method can find clusters not spherically shaped.

4 Simulated Data Set

To show FPDcoC performance it has been applied on a simulated data set. Each block has been normally generated and mean vector c_k uniformly generated. In each block $x_i \sim N_k(c_k, I)$, where I is a $J \times J$ identity matrix, rows and columns are randomly sorted. The final data set has $n = 350$ and $J = 10$. Figure 1a shows a heat-map of the original data set; the white corresponds to the minimum value of the data set, the black to the maximum value. Figure 1b represents co-clustering results on the simulated data set on 50 replications. Units and variables are sorted according to their cluster membership. In this example can be clearly seen the block structure of the data set made by height blocks. The advantage of simulated data set is that the block membership is a priori known, it allows us to measure the quality of the results using the error rate (ratio between the number of misclassified and the total number of units/variables). On the simulated data set the error rate is 0 that corresponds to a perfect classification.

5 Management Commentaries Data Set

The management commentaries data set is the result of a survey on listed companies in the Italian stock exchange. The management commentary is an important part of the official budget of listed companies, often it is not used in analysis because it is a nonstructured text but it may contain import information. Among the 406 Italian listed companies in 2009, Spano and Triunfo [18] select 25 companies and processed the management commentaries of the selected companies using ad hoc techniques. The matrix obtained is a frequency matrix (25 companies on rows, 81 words on columns). Results are shown in Fig. 2.

The resulting clustering structure is made by four clusters of firms and five clusters of words, for a total of 20 blocks. Words cluster can be interpreted according to the meaning of each word and labeled as in the following list. (1) Variations and changes, (2) economic performances, (3) costs and currencies, (4) marketing, results and previsions, and (5) sectorial words. Words concerning economic performances are the most used, they can be defined as mandatory word because in the management commentaries firms are obliged to describe their economic situation. As expected, sectorial words are the less used, each firm uses the word linked to its own sector.

Looking at cluster of firms it can be noticed that firms linked to energetic sector are not in the same cluster with the exception of the firms of gas sector that are all in cluster 4. The first cluster of firms is mainly made by editorial and communication sector firms. They use many words linked to variations and changes and they also describe the cost structure, their results, marketing strategies and some previsions. These firms have a rich language, they use all groups of words. The second cluster is the smallest, only five firms, this cluster contains only big firms, mainly of energetic sector only one, Mediaset, is linked with telecommunications. These firms have the richest language, they use all group of words with the highest frequency, they go into deep of their economic situation. The firms concerning products belong to

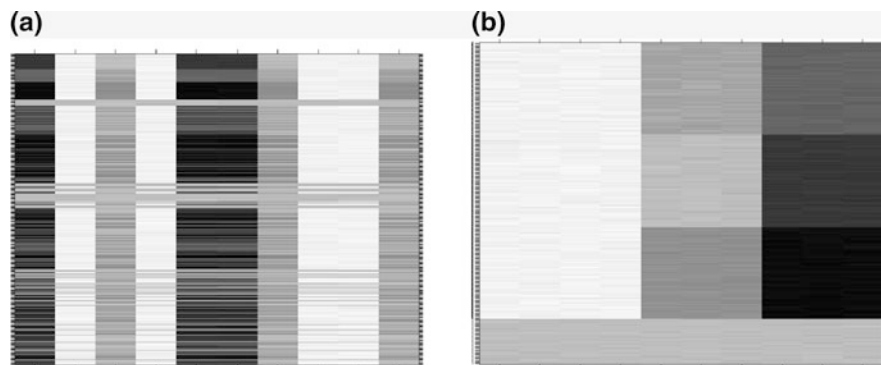


Fig. 1 Heat map of the simulated data set. (a) Original data set and (b) data set sorted according to FPDcoC results

cluster 3. These firms have the poorest language, they use few words with a very low frequency. Cluster 4 groups all firms producing gas with same firms of services (with the exception of four firms). Group 4 mainly uses words concerning economic performance and changes.

6 Conclusion

In this paper a new factor co-clustering method, FPDcoC, is proposed. The method is based on factor PD-clustering, a clustering technique that finds a factor transformation of the data and a clustering structure on the transformed data optimizing the same criterion. FPDcoC is able to catch homogeneous blocks in the data set, it has the advantage that there is an interdependence between clusters of rows and columns. The methods convergence has been empirically demonstrated. The advantages of the method are that its performance is not affected by outliers and that it can find clusters of not spherically shaped. Furthermore, the factor step allows one to work with data set of large dimensions.

References

1. Balbi, S., Miele, R., Scepi, G.: Clustering of documents from a two-way viewpoint. In: 10th International Conference on Statistical Analysis of Textual Data (2010)
2. Bock, H.H.: Two-way clustering for contingency tables: maximizing a dependence measure. In: Schader, M., et al. (eds.) *Between Data Science and Applied Data Analysis. Studies in Classification, Data Analysis, and Knowledge Organization*. Springer, Heidelberg (2003)
3. Charrad, M.M., Ahmed, M.B.: Simultaneous clustering: a survey. In: *Pattern Recognition and Machine Intelligence*, pp. 370–3775. Springer, Berlin, Heidelberg (2011)
4. Dhillon, I.S.: Co-clustering documents and words using bipartite spectral graph partitioning. In: *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 269–274 (2001)
5. Dhillon, I.S., Mallela, S., Modha, D.S.: Information-theoretic co-clustering. In: *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 89–98 (2003)
6. Gaul, W., Schader, M.: A new algorithm for two-mode clustering. In: Bock, H.H., Polasek, W. (eds.) *Data Analysis and Information Systems*, pp. 15–23. Springer, Heidelberg (1996)
7. Getz, G., Levine, E., Domany, E.: Coupled two-way clustering analysis of gene microarray data. *Proc. Natl. Acad. Sci.* **97**(22), 12079–12084 (2000)
8. Getz, G., Gal, H., Kela, I., Nottelman, D.A., Domany, E.: Coupled two-way clustering analysis of breast cancer and colon cancer gene expression data. *Bioinformatics* **19**(9), 1079–1089 (2003)
9. Govaert, G.: Simultaneous clustering of rows and columns. *Control Cybern.* **24**(4), 437–458 (1995)
10. Greenacre, M.J.: Clustering the rows and columns of a contingency table. *J. Classif.* **5**, 39–51 (1988)
11. Hartigan, J.A.: Direct clustering of a data matrix. *J. Am. Stat. Assoc.* **5**, 123–129 (1972)
12. Iyigun, C.: Probabilistic Distance Clustering. Ph.D. thesis, New Brunswick Rutgers, The State University of New Jersey (2007)
13. Krolak-Schwerdt, S.: Two-mode clustering methods: compare and contrast. In: Vichi, M., Schader, M., Gaul, W. (eds.) *Between Data Science and Applied Data Analysis: Studies in*

- Classification, Data Analysis, and Knowledge Organization, pp. 270–278. Springer, Heidelberg (2003)
14. Kroonenberg, P.M.: Applied Multiway Data Analysis. Ebooks Corporation, Hoboken, NJ (2008)
 15. Mirkin, B., Arabie, P., Hubert, L.J.: Additive two- mode clustering: the error-variance approach revisited. *J. Classif.* **12**, 243–263 (1995)
 16. Pollard, K.S., Van der Laan, M.J.: Statistical inference for simultaneous clustering of gene expression data. *Math. Biosci.* **176**(1), 99–121 (2002)
 17. Rocci, R., Vichi, M.: Two-mode multi-partitioning. *Comput. Stat. Data Anal.* **52**(4), 1984–2003 (2008)
 18. Spano, M., Triunfo, N.: La relazione sulla gestione delle società italiane quotate sul mercato regolamentare. In: *JADT 2012: 11es Journées internationales d'Analyse statistique des Données Textuelles* (2012)
 19. Tang, C., Zhang, L., Zhang, A., Ramanathan, M.: Interrelated two-way clustering: an unsupervised approach for gene expression data analysis. In: *Bioinformatics and Bioengineering Conference, 2001. Proceedings of the IEEE 2nd International Symposium on IEEE*, pp. 41–48 (2001)
 20. Tortora, C.: Non-hierarchical clustering methods on factorial subspaces. Ph.D. thesis, Università di Napoli Federico II (2011)
 21. Tortora, C., Marino, M.: Robustness and stability analysis of factor PD-clustering on large social datasets. In: Vicari, D., Okada, A., Ragozini, G., Weihs, C. (eds.) *Analysis and Modeling of Complex Data in Behavioral and Social Sciences*, pp. 273–281. Springer, New York (2014)
 22. Tortora, C., Gettler Summa, M., Palumbo, F.: Factor PD-clustering. In: Alfred, U., Berthold, L., Dirk, V. (eds.) *Algorithms from and for Nature and Life*, pp. 115–123. (2013)
 23. Van Mechelen, I., Bock, H.H., De Boeck, P.: Two-mode clustering methods: a structured overview. *Stat. Methods Med. Res.* **13**(5), 363–394 (2004)
 24. Vichi, M.: Double k-means clustering for simultaneous classification of objects and variables. In: Borra, S., Rocci, R., Vichi, M., Schader, M. (eds.) *Advances in Classification and Data Analysis*, pp. 43–52. Springer, Heidelberg (2001)
 25. Vichi, M., Kiers, H.A.L.: Factorial k-means analysis for two way data. *Comput. Stat. Data Anal.* **37**, 29–64 (2001)

Part III

Sampling and Estimation Methods

M-Quantile Small Area Estimation for Panel Data

Annamaria Bianchi

Abstract

Economic indicators need to be estimated at regional level. Small area estimation based on M-quantile regression has recently been introduced by Chambers and Tzavidis (*Biometrika* 93:255–268, 2006) and it has proved to provide a valid alternative to traditional methods. Thus far, this method has only been applied to cross-sectional data. However, it is well known that the use of panel data may provide significant gains in terms of efficiency of the estimators. This paper explores possible extensions of M-quantile-based small area estimators to the panel data context. A model-based simulation study is presented.

1 Introduction

In recent years, there has been an increasing demand by policy makers for estimates of economic indicators at regional level. Unfortunately, limited founding resources for the design of sample surveys often lead to small sample sizes within these domains. As a consequence, direct estimators (which use only data from sample units in the domain) cannot be applied since they yield estimates with unacceptable standard errors. These problems are typically overcome by the use of small area techniques. This is an approach based on models that borrow strength in making an estimate for one small area from sample survey data collected in other small areas and/or at other time periods. The most popular class of models for small area estimation is based on random effects models, which include random area effects to account for between area variation.

A. Bianchi (✉)

Università degli Studi di Bergamo, via dei Caniana 2, 24127 Bergamo, Italy
e-mail: annamaria.bianchi@unibg.it

National Statistical Institutes periodically collect survey data using a panel approach, that is, multidimensional data involving measurements over time (phenomena are observed over multiple time periods for the same units). An important example of (rotating) panel is the Labour Force Survey. Panel data combine the individual dimension with the time dimension, thereby augmenting the information of the data with respect to a cross-section approach. For this reason panel data analysis presents many benefits. It allows to control for individual (and time) unobserved heterogeneity, and hence allows to isolate the longitudinal variability of the investigated phenomena from the variability due to the different characteristics of the responding units. Moreover, panel data are more informative since there is more variability and the estimates are therefore more efficient. Starting from this situation and given the need of regional indicators, in this paper we focus on small area models that borrow strength across both small areas and times.

In the small area context, it is well known that for such panel surveys considerable gains in efficiency can be achieved by borrowing strength across both small areas and times. Thus far, the use of longitudinal data for purposes of small area estimation is concentrated mostly on the area level models [7]. The possible reason for this is that in many countries, and especially in the USA, the infrastructure of the official statistics does not support longitudinal data sets at individual level. On the other hand, research on small area estimation from unit-level panel data is clearly needed, because aggregating individual level data to adapt for area level models may cause unnecessary loss of information. At the unit level, an appropriate model for panel data must take the covariance of the repeated observations from the same unit into account. One simple model that can be adapted to this purpose is the two-fold nested error regression model proposed by Stukel and Rao, 1999 [9]. The small area means then are estimated by the empirical best linear unbiased predictor (EPLUP). Refer to Rao [6] for more details.

Recently, an alternative unit-level approach to small area estimation based on M-quantile (MQ) regression has been proposed by Chambers and Tzavidis [3]. The advantages of MQ-based small area estimators with respect to traditional random effects models are that they do not depend on strong distributional assumptions and that they are outlier robust. The initial estimator proposed in Chambers and Tzavidis [3] has subsequently been extended in various ways (see [8] and [10]). However, to the best of our knowledge, up to now this technique has only been applied to cross-sectional data. The gains in efficiency that can be obtained using panel data have not been explored yet.

The aim of this research is to explore whether MQ small area estimation methods based on panel data provide significant improvements over the more classical MQ technique for cross-sectional data. The longitudinal extension of MQ-based small area estimation is expected to be useful in practical applications, especially in the context of income studies. Recently, MQ methods have been intensively applied to study income indicators, using especially data coming from the European Survey on Income and Living Conditions (EU-SILC) waves [4]. EU-SILC is a panel survey. The proposed method could provide an improvement of small area MQ-based estimates for this kind of data.

The rest of the paper is organized as follows. In Sect. 2 a review of MQ small area estimation for cross-sectional data is presented. Section 3 provides an extension to panel data. The performance of the proposed method is analyzed in Sect. 4 by means of a model-based simulation study. Section 5 concludes.

2 M-Quantile Small Area Estimation for Cross-Sectional Data

Suppose that a population U of size N is divided into non-overlapping domains of size $N_j, j = 1, \dots, d$. Assume that a sample s is available and denote n_j the sample size in area j and s_j (r_j) the sampled (non-sampled) population units in the area. Let y_{ij} denote the value of the variable of interest for unit belonging to the small area j ($j = 1, \dots, d, i = 1, \dots, n_j$). Assume that the values of y are available for each unit of the sample and that a vector of auxiliary variables \mathbf{x}_{ij} is available for each unit of the population. We are interested in predicting small area means for the target variable for each small area: $m_j = N_j^{-1} \sum_{i \in s_j \cup r_j} y_{ij}$ ($j = 1, \dots, d$).

A recently introduced approach to small area estimation is based on MQ regression. MQ regression [2] provides a “quantile-like” generalization of regression based on influence functions. For fixed q ($0 < q < 1$), the linear MQ model of order q of the conditional distribution of y given x is given by

$$MQ_q(y_{ij}|\mathbf{x}_{ij}) = \mathbf{x}_{ij}^T \boldsymbol{\beta}_q.$$

An estimate $\hat{\boldsymbol{\beta}}_q$ of $\boldsymbol{\beta}_q$ is obtained by solving the following equations (in $\boldsymbol{\beta}$)

$$\sum_{j=1}^d \sum_{i=1}^{n_j} \psi_q \left(\frac{y_{ij} - \mathbf{x}_{ij}^T \boldsymbol{\beta}}{s_n} \right) \mathbf{x}_{ij} = \mathbf{0},$$

where $\psi_q(u) = |q - I(u < 0)|\psi(u)$, ψ is an appropriately chosen influence function (such as the Huber Proposal 2), and s_n is a robust estimate of scale such as the median absolute deviation. The idea underlying MQ-based small area estimation is the following. The conditional variability across the population can be characterized by the so-called MQ coefficients of the population units. For unit i in small area j with values $(\mathbf{x}_{ij}, y_{ij})$ this coefficient is defined as the value q_{ij} such that $MQ_{q_{ij}}(y_{ij}|\mathbf{x}_{ij}) = y_{ij}$ —that is, q_{ij} is the order of the MQ passing through the point $(\mathbf{x}_{ij}, y_{ij})$. If a hierarchical structure does explain part of the variability in the population, then it is expected that units belonging to the same area have similar coefficients. It is therefore natural to characterize each small area j by means of an indicator θ_j defined here as the mean of the population MQ coefficients belonging to that area $\theta_j = N_j^{-1} \sum_{i \in s_j \cup r_j} q_{ij}$. This coefficient identifies an MQ regression plane $MQ_{\theta_j}(y_{ij}|\mathbf{x}_{ij})$ characteristic of that area, which allows to predict unobserved data in the area. Such predicted values are then used to construct estimates of m_j . When $\boldsymbol{\beta}_q$ is a sufficiently smooth function of q , the naive estimator has first been proposed by

Chambers and Tzavidis [3]

$$\hat{m}_j = N_j^{-1} \left[\sum_{i \in s_j} y_{ij} + \sum_{i \in r_j} \mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}}_{\hat{\theta}_j} \right]. \quad (1)$$

Unfortunately estimator (1) turned out to be biased, especially when there are large outlying values in the dataset. For this reason Tzavidis et al. [10] proposed a bias-adjusted estimator

$$\hat{m}_j^{\text{CD}} = N_j^{-1} \left[\sum_{i \in s_j} y_{ij} + \sum_{i \in r_j} \mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}}_{\hat{\theta}_j} + \frac{N_j - n_j}{n_j} \sum_{i \in s_j} \{y_{ij} - \mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}}_{\hat{\theta}_j}\} \right]. \quad (2)$$

Refer to Tzavidis et al. [10] and Salvati et al. [8] for other possible estimators, together with the corresponding MSE estimators.

3 M-Quantile Small Area Estimation for Panel Data

Let now y_{ijt} denote the value of the variable of interest for unit i belonging to the small area j at time t ($j = 1, \dots, d$, $t = 1, \dots, T$, $i = 1, \dots, n_{jt}$). Denote \mathbf{x}_{ijt} the corresponding covariates known at the population level. We are now interested in predicting small area means for the target variable at the final time T : $m_{jT} = N_j^{-1} \sum_{i \in s_j \cup r_j} y_{ijT}$ ($j = 1, \dots, d$). In order to extend the MQ-based small area technique to panel data, the first step is to extend MQ regression to panel data. For a given q , the simplest MQ panel data model is defined by

$$\text{MQ}_q(y_{ijt} | \mathbf{x}_{ijt}) = \mathbf{x}_{ijt}^T \boldsymbol{\beta}_q, \quad t = 1, \dots, T. \quad (3)$$

It is worth making a few comments on the model. First, cross-sectional MQ is a special case of (3). Indeed, it is enough to consider $T = 1$. Second, model (3) may appear quite restrictive since $\boldsymbol{\beta}_q$ is the same in each time period. However, by appropriately interacting covariates with time dummies, parameters are allowed to change over time. Similarly, one can include different parameters for different groups. Third, this model is based on the assumption of homogeneity of the coefficients across individuals (no individual effects). Of course, this assumption should be tested when applying the model to real data. Unfortunately, hypothesis testing theory for MQ regression has not been developed yet, even in the cross-sectional context. For this reason, the poolability assumption cannot be tested in the present paper, as it would require the development of hypothesis testing theory for the cross-sectional case first. The author is currently working on this topic. It is expected that similar tests to the ones used in OLS regression (though adapted to the MQ context) could be used. Fourth, with the introduced simplification, the panel data model covered here, while having many useful applications, does not fully exploit the replicability over time. A further step in the improvement of the estimates

is expected by allowing explicitly for unobserved individual effects. Unfortunately, this is a complicated scenario in the case of MQ regression. The same problem is faced by quantile regression for panel data and it has been treated recently. Refer to Wooldridge [11] and references therein.

The natural estimator for β_q is the pooled MQ regression estimator $\hat{\beta}_q$, which solves

$$\sum_{j=1}^d \sum_{t=1}^T \sum_{i=1}^{n_{jt}} \psi_q \left(\frac{y_{ijt} - \mathbf{x}_{ijt}^T \beta}{s_n} \right) \mathbf{x}_{ijt} = \mathbf{0}.$$

Under the previous assumption, this kind of regression allows to look at the dynamic relationship and is expected to increase the efficiency of estimators. Asymptotic theory for $\hat{\beta}_q$ (for cross-sectional dimension tending to infinity and fixed time dimension) follows from cross-sectional MQ regression asymptotics. Refer to Bianchi and Salvati [1]. The estimates of the small area means are then computed by applying estimators (1) and (2) using β estimates from the panel MQ regression and data from the final time T . The idea is that the improved β estimates should imply better estimates for the small area means too.

4 Simulation Study

We carry out a Monte Carlo simulation study to investigate the performance of the introduced estimators for panel data. The task is to estimate the area mean of the response variable at the last (“current”) time point in the panel.

The simulation is model-based in the sense that in each replication we generated a different longitudinal population of $d = 15$ areas and $T = 5$ time periods using a specified model. For each replication we drew a sample of $n = 250$ units. The procedure was repeated $K = 1000$ times. No drop-outs or non-response is introduced. The sampling method is the simple random sampling. Hence the regional sample sizes are not controlled in any way and they simply reflect the regional population sizes.

The model equation used in generating the population y values for each replicate was

$$y_{ijt} = \mu_t + x_{ijt} + u_i + v_{ij} + e_{ijt},$$

where

- \mathbf{x} values within area j were drawn independently at each replication as $\mathbf{x} = [x_{ij1}, x_{ij2}, x_{ij3}, x_{ij4}, x_{ij5}] \sim N(\mu = \mu_j, \sigma^2 = \mu_j^2/36, \rho = 0.5)$,
- the population mean μ_j of the \mathbf{x} values in area j was chosen at random without replacement from the integers between 40 and 200 and was held fixed over all simulations

- the random area effects u_i , the random individual effects v_{ij} , and the error term e_{ijt} were independently generated from normal distributions: $N(0, \sigma_u^2)$, $N(0, \sigma_v^2)$, and $N(0, \sigma_e^2)$, respectively
- the time contribution was considered fixed $\mu_t = [1/6, 3/6, 5/6, 7/6, 9/6]$.

The factors whose effect on the estimation performance are evaluated and which define the design of the simulation experiment are the intra-area correlation and the intra-unit correlation. The intra-area correlation of the units measures homogeneity of units belonging to the same area. For two different units ij and ij' ($j \neq j'$) from the same area i , it is defined as

$$\frac{\sigma_u^2}{\sigma_u^2 + \sigma_v^2 + \sigma_e^2}.$$

The intra-unit correlation measures the level of correlation of repeated observations from the same unit. For two observations y_{ijt} and $y_{ij't'}$ from the same unit at different time points $t \neq t'$, it is defined as

$$\frac{\sigma_u^2 + \sigma_v^2}{\sigma_u^2 + \sigma_v^2 + \sigma_e^2}.$$

The different covariance structures of the generated populations I–IV were determined by different variances for the random effects. The variances and the resulting correlations are shown in Table 1.

Notice that two levels are considered both for the intra-area correlation (high in Populations I and II and low in Populations III and IV) and the intra-unit correlation (high in Populations I and III and low in Populations II and IV).

In each replication, the estimates are computed by fitting a panel MQ model to the data

$$y_{ijt} = \alpha_t + \beta x_{ijt} + e_{ijt}$$

and then applying the estimators (1) and (2). For comparison purposes, we also consider the traditional MQ estimator for cross-sectional data, that is, we fit a

Table 1 Variances of the random effects and corresponding intra-area and intra-unit correlations in the generated populations (I–IV)

	I	II	III	IV
Area variance	3	3	1	1
Individual variance	4	1	5	1
Error variance	1	4	2	6
Intra-area corr.	0.375	0.375	0.125	0.125
Intra-unit corr.	0.875	0.5	0.75	0.25

cross-sectional MQ model (which only uses the cross-sectional data of the last time point T)

$$y_{ijT} = \alpha_t + \beta x_{ijT} + e_{ijT}.$$

The performance of the proposed estimators for the small area means are examined by computing for each small area the relative bias (RB) and the relative root mean squared error (RRMSE). Denoting \hat{m}_{ki} the estimate of m_i obtained from the k -th replication, they are computed as follows

$$RB_i = \frac{\sum_{k=1}^K (\hat{m}_{ki} - m_{ki})/K}{\sum_{k=1}^K m_{ki}/K} \times 100$$

where k indicates the iteration number. The root mean squared error is computed as

$$RRMSE_i = \frac{\sqrt{\sum_{k=1}^K (\hat{m}_{ki} - m_{ki})^2 / K}}{\sum_{k=1}^K m_{ki} / K} \times 100.$$

Table 2 shows the mean and the five-point summary (minimum, first quartile, median, third quartile, and maximum) of the distribution of the RB and RRMSE over the small areas.

From the results we see that panel MQ improves the estimates in the case of the naive estimator. The highest improvements are observed for Populations II and IV, that is, when the intra-unit correlation is low. This situation corresponds to high heterogeneity of observations from the same individual, yielding more informative data. For the CD estimator only slight improvements are observed using the panel MQ model. Further investigations are needed to assess the effect of the length of the panel and the specification of the model.

5 Concluding Remarks

This paper can be viewed as a first exploration of M-quantile-based small area estimators in the panel data context. The contribution of this study is twofold. First, a simple extension of MQ regression to panel data is proposed. Second, this longitudinal MQ models are used in the small area estimation. As regards the first contribution, the extension is limited to a simple model, not including unobserved individual effects. Indeed, the inclusion of unobserved individual effects is rather complicated in the MQ context, leading to the incidental parameters problem. The same problem is faced by quantile regression for panel data and it has been treated recently [11]. Moreover, hypothesis testing theory has not been developed for cross-sectional MQ regression yet. This prevents the development of proper testing on the presence of individual effects as well as general model specification. Further research is needed in order to provide more extensive methodological developments.

Table 2 Relative bias (RB) and relative root mean squared error (RRMSE) in model-based simulation study

	NAIVE.PANEL		NAIVE.CS		CD.PANEL		CD.CS	
	RB	RRMSE	RB	RRMSE	RB	RRMSE	RB	RRMSE
<i>Population I</i>								
Min.	-0.0194	0.1978	-0.0487	0.2897	-0.0621	0.3079	-0.0610	0.3086
1st Qu.	-0.0107	0.3383	-0.0087	0.5303	-0.0268	0.6233	-0.0272	0.6240
Median	0.0023	0.5217	0.0080	0.7160	0.0114	0.8498	0.0114	0.8508
Mean	0.0091	0.4800	0.0148	0.7115	0.0116	0.8129	0.0118	0.8134
3rd Qu.	0.0134	0.5901	0.0271	0.9232	0.0372	1.0570	0.0372	1.0570
Max.	0.0861	0.7415	0.1236	1.1980	0.1317	1.3160	0.1333	1.3160
<i>Population II</i>								
Min.	-0.0810	0.3394	-0.0989	0.3902	-0.0628	0.3012	-0.0623	0.3012
1st Qu.	-0.0096	0.3856	0.0022	0.4558	-0.0010	0.3967	-0.0017	0.3959
Median	0.0109	0.4630	0.0136	0.5409	0.0054	0.4739	0.0071	0.4744
Mean	0.0112	0.7318	0.0233	0.9638	0.0259	0.9141	0.0263	0.9138
3rd Qu.	0.0210	0.8543	0.0350	1.0750	0.0219	1.0130	0.0224	1.0110
Max.	0.1414	2.3700	0.2773	3.7930	0.2463	4.0610	0.2477	4.0630
<i>Population III</i>								
Min.	-0.0699	0.2673	-0.0992	0.3104	-0.1124	0.2695	-0.1153	0.2703
1st Qu.	-0.0256	0.4064	-0.0373	0.4940	-0.0220	0.4650	-0.0222	0.4649
Median	-0.0102	0.5273	-0.0058	0.6304	-0.0056	0.5512	-0.0058	0.5500
Mean	-0.0063	0.6203	-0.0087	0.7716	-0.0107	0.7053	-0.0110	0.7055
3rd Qu.	0.0114	0.8243	0.0122	1.0030	0.0191	0.8593	0.0192	0.8594
Max.	0.0913	1.2400	0.1149	1.5920	0.0288	1.6820	0.0282	1.6830
<i>Population IV</i>								
Min.	-0.0446	0.3108	-0.0824	0.3954	-0.1329	0.3901	-0.1329	0.3897
1st Qu.	-0.0029	0.3954	-0.0052	0.5437	-0.0137	0.6282	-0.0138	0.6275
Median	0.0175	0.5256	0.0201	0.7852	0.0113	0.8171	0.0115	0.8164
Mean	0.0222	0.5791	0.0189	0.8642	0.0081	0.9863	0.0072	0.9863
3rd Qu.	0.0511	0.7324	0.0466	1.1090	0.0400	1.2310	0.0392	1.2350
Max.	0.0821	1.1100	0.1264	1.8230	0.0933	2.1630	0.0871	2.1600

Nevertheless, the results from the simulation study are rather promising. Moreover, the recent application of small area MQ-based estimates to EU-SILC panel data can make the application of the proposed extension interesting. The method could be extended to other economic panel data as well. For example, it could be an alternative to other methods in the context of the Labour Force Survey (see, e.g., [5]).

Acknowledgements Paper supported by the ex 60% University of Bergamo, Biffignandi grant, and PAADEL project (Lombardy Region-University of Bergamo joint project). The author is thankful to the anonymous referee for providing very useful comments.

References

1. Bianchi, A., Salvati, N.: Asymptotic properties and variance estimators of the M-quantile regression coefficients estimators. *Commun. Stat. Theory Methods* **44**, 2416–2429 (2015)
2. Breckling, J., Chambers, R.: M-quantiles. *Biometrika* **75**, 761–771 (1988)
3. Chambers, R., Tzavidis, N.: M-quantile models for small area estimation. *Biometrika* **93**, 255–268 (2006)
4. Fabrizi, E., Giusti, C., Salvati, N., Tzavidis, N.: Mapping average equivalized income using robust small area methods. *Pap. Reg. Sci.* **93**, 685–701 (2014)
5. Falorsi, P.D., Falorsi, S., Russo, A.: Small area estimation at provincial level in the Italian Labour Force Survey. *J. Ital. Stat. Soc.* **11**, 93–109 (1998)
6. Rao, J.N.K.: *Small Area Estimation*. Wiley, New York (2003)
7. Rao, J.N.K., Yu, M.: Small area estimation by combining time series and cross-sectional data. *Can. J. Stat.* **22**, 511–528 (1994)
8. Salvati, N., Tzavidis, N., Pratesi, M., Chambers, R.: Small area estimation via M-quantile geographically weighted regression. *Test* **21**, 1–28 (2012)
9. Stukel, D.M., Rao, J.N.K.: On small-area estimation under two-fold nested error regression models. *J. Stat. Plan. Inference* **78**, 131–147 (1999)
10. Tzavidis, N., Marchetti, S., Chambers, R.: Robust estimation of small area means and quantiles. *Aust. N. Z. J. Stat.* **52**, 167–186 (2010)
11. Wooldridge, J.: *Econometric Analysis of Cross Section and Panel Data*. MIT Press, Cambridge (2010)

A Two-Part Geoadditive Small Area Model for Geographical Domain Estimation

Chiara Bocci, Alessandra Petrucci, and Emilia Rocco

Abstract

We are interested in estimating small domain means of a response variable that shows a spatial trend and has a continuous skewed distribution with a large number of values clustered at zero. This kind of variable can occur in many surveys, like business or agricultural surveys: examples are the quantity of crops produced or the amount of land allocated for their production collected by the Farm Structure Survey driven by the Italian Statistical Institute. The small sample size within the areas requires the use of small area model dependent methods to increase the effective area sample size by using census and administrative auxiliary data. To account simultaneously for the excess of zeros, the skewness of the distribution and the possible spatial trend of the data, we present a two-part geoadditive small area model. An application to the estimation of the per-farm average grapevine production in Tuscany at Agrarian Region level shows the satisfactory performance of the model.

C. Bocci

IRPET - Regional Institute for Economic Planning of Tuscany, Via Pietro Dazzi 1, 50141 Firenze, Italy

e-mail: chiara.bocci@irpet.it

A. Petrucci • E. Rocco (✉)

Department of Statistics, Computer Science, Applications “G. Parenti”, University of Florence, Viale Morgagni 59, 50134 Firenze, Italy

e-mail: alessandra.petrucci@unifi.it; rocco@disia.unifi.it

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*, Studies in Theoretical and Applied Statistics,

DOI 10.1007/978-3-319-27274-0_12

1 Introduction

Since year 2000, the Italian Agricultural Census driven by the Italian Statistical Institute (ISTAT) has given the opportunity to georeference the farms introducing a new challenge for statistical analysis of phenomena concerning the agricultural field. The geographical location of each farm can constitute a particularly useful information for the analysis of many phenomena concerning the agricultural field [2]. In fact almost all the variables of interest in agricultural field have a spatial trend and estimates for particular geographical areas are often required by policy makers to support fund allocation, regional planning, etc. This article, motivated by the real problem of providing reliable mean estimates for small geographical domains of variables concerning agricultural phenomena, shows how the geographical location of the farms can be a very important auxiliary variable in a model based approach to the small area estimation (SAE) problem.

SAE methods are concerned with producing reliable estimates of characteristics of interest such as means, counts, etc., for areas or domains for which only small samples or no samples are available. Sample surveys are usually planned to produce estimates for larger domains or areas and, due to cost and operational considerations it is usually not possible to procure a large enough overall sample size to support direct estimates of adequate precision for all areas of interest. This makes it necessary to borrow information across related areas through indirect estimation based on models, using auxiliary information such as census and administrative data. The most popular class of SAE models is the linear mixed models that include independent random area effects to account for between area variation beyond that explained by auxiliary variables [1, 7]. Numerous extensions to this set-up have been presented in the literature to include various generalized linear models and/or more complicated random effects structures.

Our aim is to build a SAE model for estimating the mean of a variable with the following characteristics: (a) a large peak at zero; (b) a skewed distribution of the non-zero values; and (c) related to georeferenced units. Each one of these characteristics of the data has been separately addressed in the literature, in the context of SAE or in other fields of research, and by merging these proposals we want to define a SAE model able to account simultaneously for all these aspects of the data. In the following, we give a brief excursus of this literature.

First, the excess zeros in data are usually described by the zero-inflated (ZI) regression models that mix a degenerate distribution with point mass of one at 0 with a simple regression model based on a standard distribution. This is realized considering a pair of regression models: a model, usually logit or probit, for the probability of non-zero response and a conditional regression model for the mean response given that it is non-zero. The ZI models have been originally developed to analyze count data, but they are also extended to situation in which a huge number of zeros occur in continuous data [8, 12]. Frequently, in this context they are known as two-part models. In the context of SAE methods Pfeiffermann et al. [14] describe the problem of zero-inflated data (not skewed) under a two-part random effects model

using a Bayesian approach. Chandra and Sud [5] consider the same framework but adopt a frequentist approach.

Second, when the data are skewed the relationship between the response variable and the auxiliary variable may not be linear in the original scale, but can be linear in a transformed scale, e.g., the logarithm scale. In such case, SAE methods based on linear models produce inefficient estimates and, using the log-transformed models, there are alternative approaches to obtain better indirect predictors for small area mean [3, 16]. Zero-inflated log-normal models have been largely applied for the analysis of longitudinal data [8, 12] and recently their use in the context of SAE has been suggested by Chandra and Chambers [4].

Last, if the data are characterized by a non stationary spatial trend, the use of “global” dependence models, that assume independency of the data from their spatial location, can generate spatially autocorrelated residuals and bring often to wrong conclusions. An adequate use of the geographic information in more complex models which take into account the spatial variability can help to understand the underlying phenomenon. In literature several models have been proposed to simultaneously incorporate the spatial distribution of the study variable and the other covariate effects. Geoadditive models [10], in particular, merge an additive model [9], that accounts for the relationship between the variables, and a Kriging model [6], that accounts for the spatial distribution, under the linear mixed model framework. The mixed model structure allows to easily include the small area specific effect as an additional random components, obtaining a geoadditive SAE model [13].

Merging these approaches, we derive the two-part geoadditive small area model illustrated in Sect. 2. The model is fitted to data on Tuscany grapevine production collected by the Farm Structure Survey (FSS) carried out by ISTAT in 2003 and the results are presented in Sect. 3. Section 4 concludes with some final remarks.

2 Model and Small Area Mean Predictor

2.1 Two-Part Geoadditive Small Area Model

Let y_{ij} denote a non-negative semi-continuous skewed response variable for unit j ($j = 1, \dots, N_i$) in small area i ($i = 1, \dots, m$), with $\sum_{i=1}^m N_i = N$, placed at spatial location \mathbf{s}_{ij} ($\mathbf{s} \in R^2$). This response variable can be recoded as the product $y_{ij} = I_{ij}y'_{ij}$ of two variables

$$I_{ij} = \begin{cases} 1 & \text{if } y_{ij} > 0 \\ 0 & \text{if } y_{ij} = 0 \end{cases} \quad \text{and} \quad y'_{ij} = \begin{cases} y_{ij} & \text{if } y_{ij} > 0 \\ \text{irrelevant} & \text{if } y_{ij} = 0 \end{cases} .$$

We model these variables by a pair of uncorrelated random effects models, one for the logit of probability $\pi_{ij} = P(I_{ij} = 1) = P(y_{ij} > 0)$ and one for the mean conditional response $E[y'_{ij}|I_{ij} = 1]$. Both models include individual covariates, as well as area random effects that account for variations not explained by the

covariates. Moreover, as we assume that y may present a relevant spatial pattern, both models include also a smooth, non-parametrically specified spatial trend.

The logit model is

$$\eta_{ij} = \log \frac{\pi_{ij}}{1 - \pi_{ij}} = \beta_{0t} + \mathbf{t}_{ij}^T \boldsymbol{\beta}_t + h(\mathbf{s}_{ij}) + u_i \quad u_i \sim N(0, \sigma_u^2) \quad (1)$$

where \mathbf{t}_{ij} is a vector of p linear covariates, $h(\cdot)$ is an unspecified bivariate smooth function depending on geographical unit coordinates \mathbf{s}_{ij} , and u_i is a area-specific random effect associated with area i . Representing $h(\cdot)$ with a low rank thin plate spline [15] with K knots $(\kappa_1, \dots, \kappa_K)$, that is,

$$h(\mathbf{s}_{ij}) = \beta_{0s} + \mathbf{s}_{ij}^T \boldsymbol{\beta}_s + \sum_{k=1}^K \gamma_k b(\mathbf{s}_{ij}, \kappa_k), \quad (2)$$

model (1) can be written as a mixed model [10, 13]:

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{D}\mathbf{u} \quad (3)$$

where:

- \mathbf{X} is the $N \times (p + 3)$ fixed effects matrix with rows $[1, \mathbf{t}_{ij}^T, \mathbf{s}_{ij}^T]$;
- $\mathbf{Z} = [C(\mathbf{s}_{ij} - \boldsymbol{\kappa}_k)]_{\substack{1 \leq j \leq N_i \\ 1 \leq i \leq m \\ 1 \leq k \leq K}} [C(\boldsymbol{\kappa}_h - \boldsymbol{\kappa}_k)]_{\substack{1 \leq h \leq K \\ 1 \leq k \leq K}}^{-1/2}$ with $C(\mathbf{r}) = \|\mathbf{r}\|^2 \log \|\mathbf{r}\|$ is the matrix of the thin plate spline basis functions $b(\mathbf{s}_{ij}, \kappa_k)$;
- \mathbf{D} is the $N \times m$ area-specific random effects matrix with rows \mathbf{d}_{ij} containing indicators taking value 1 if observation j is in area i and 0 otherwise;
- $\boldsymbol{\beta} = [\beta_{0t} + \beta_{0s}, \boldsymbol{\beta}_t^T, \boldsymbol{\beta}_s^T]^T$ is a vector of unknown coefficients;
- $\mathbf{u} \sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}_m)$ and $\boldsymbol{\gamma} \sim N(\mathbf{0}, \sigma_\gamma^2 \mathbf{I}_K)$ are respectively the vector of the m area-specific random effects and the vector of the K thin plate spline coefficients.

The mixed model for the mean conditional response is assumed linear in the log-transformed scale. That is, for units with positive response we assume:

$$l_{ij} = \log(y'_{ij}) = \beta_{0t}^* + \mathbf{t}_{ij}^{*T} \boldsymbol{\beta}_t^* + h^*(\mathbf{s}_{ij}) + u_i^* + e_{ij} \quad \begin{matrix} u_i^* \sim N(0, \sigma_{u^*}^2) \\ e_{ij} \sim N(0, \sigma_e^2) \end{matrix} \quad (4)$$

In analogy with model (1), \mathbf{t}_{ij}^* represents individual covariates, u_i^* denotes the random area effect, and $h^*(\cdot)$ is an unspecified bivariate smooth function depending

on geographical unit coordinates s_{ij} . Note that the covariates t_{ij} in Eq. (1) and the covariates t_{ij}^* in Eq. (4) may be equal or may differ totally or partially. Representing $h^*(\cdot)$ as in (2), model (4) can be written in the mixed model form:

$$\log(y') = \mathbf{X}^* \boldsymbol{\beta}^* + \mathbf{Z}^* \boldsymbol{\gamma}^* + \mathbf{D}^* \mathbf{u}^* + \mathbf{e} \tag{5}$$

where y' is the vector of length N^* containing the response values only for the units for which $I_{ij} = 1$ and, same as before, $\boldsymbol{\gamma}^*$ and \mathbf{u}^* are respectively the vector of spline coefficients and the vector of the area-specific random effects, \mathbf{X}^* is the matrix of covariates relating to the fixed effects, \mathbf{Z}^* is the spline basis matrix, and \mathbf{D}^* is the matrix of covariates concerning the area-specific random effects.

Making the assumption of uncorrelated random effects models, the likelihood function for this mixture model factors into two terms, one for the zero and one for the non-zero data, so that it is equivalent to separately model the non-zero data y_{ij} and the indicator variable I_{ij} . Dealing with clustered data, we are aware that this assumption may not be valid since the cluster-specific random effects included into the two models may be correlated. However, in a recent paper Zhang et al. [17] compared the parameter estimates obtained adopting a two-part hierarchical model with a correlated random effects structure with those obtained fitting separately the two models and showed that they are similar. Supported by this result, in the application described in Sect. 3 we maintain the assumption that the random effects relative to the two models are uncorrelated.

2.2 Small Area Mean Predictor

We are interested in the estimation of the area means $\bar{Y}_i = N_i^{-1} \sum_{j=1}^{N_i} y_{ij}$. Notice that the means are computed over all the individuals, including individuals with zero y values, and that they may be decomposed as:

$$\bar{Y}_i = N_i^{-1} \left(\sum_{j \in S_i} y_{ij} + \sum_{j \in R_i} y_{ij} \right)$$

where S_i is the area specific sample and R_i is its complement to the area population. Under the two-part models (3) and (5) the area means predictors are defined as:

$$\hat{\bar{Y}}_i = N_i^{-1} \left(\sum_{j \in S_i} y_{ij} + \sum_{j \in R_i} \hat{y}_{ij} \right) \tag{6}$$

where \hat{y}_{ij} are the predicted values and are obtained as:

$$\hat{y}_{ij} = \frac{\exp(\mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}} + \mathbf{z}_{ij}^T \hat{\boldsymbol{\gamma}} + u_i)}{1 + \exp(\mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}} + \mathbf{z}_{ij}^T \hat{\boldsymbol{\gamma}} + u_i)} \times \hat{\lambda}_{ij}^{-1} \exp\left(\mathbf{x}_{ij}^{*T} \hat{\boldsymbol{\beta}}^* + \frac{\hat{v}_{ij}}{2}\right) \quad (7)$$

where $\hat{v}_{ij} = \text{Var}(l_{ij}) = \hat{\sigma}_e^2 + \hat{\sigma}_{u^*}^2 + \mathbf{z}_{ij}^T \hat{\sigma}_{\boldsymbol{\gamma}^*}^2 \mathbf{z}_{ij}$, and $\hat{\lambda}_{ij}$ is the bias adjustment factor for the log-back transformation suggested by Chandra and Chambers [3] and its expression is: $\hat{\lambda}_{ij} = 1 + 0.5[\hat{a}_{ij} + 0.25\hat{V}(\hat{v}_{ij})]$ where $\hat{a}_{ij} = \mathbf{x}_{ij}^T \hat{V}(\hat{\boldsymbol{\beta}}^*) \mathbf{x}_{ij}$, $\hat{V}(\hat{\boldsymbol{\beta}}^*)$ is the usual estimator of $\text{Var}(\hat{\boldsymbol{\beta}}^*)$ and $\hat{V}(\hat{v}_{ij})$ is the estimated asymptotic variance of \hat{v}_{ij} .

2.3 Mean Squared Error Estimation

Following Chandra and Sud [5], we estimate the mean squared error (MSE) of predictors (6) by adopting a parametric bootstrap approach.

For given $\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\beta}}^*, \hat{\sigma}_e^2, \hat{\sigma}_u^2, \hat{\sigma}_{u^*}^2, \hat{\sigma}_\gamma^2, \hat{\sigma}_{\gamma^*}^2$, we generate K random spline coefficients $\tilde{\boldsymbol{\gamma}}$ and $\tilde{\boldsymbol{\gamma}}^*$, m random components $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{u}}^*$, and N unit-level random errors $\tilde{\mathbf{e}}$ by drowing from the stochastic models $\boldsymbol{\gamma} \sim N(\mathbf{0}, \hat{\sigma}_\gamma^2 \mathbf{I}_K)$, $\boldsymbol{\gamma}^* \sim N(\mathbf{0}, \hat{\sigma}_{\gamma^*}^2 \mathbf{I}_K)$, $\mathbf{u} \sim N(\mathbf{0}, \hat{\sigma}_u^2 \mathbf{I}_m)$, $\mathbf{u}^* \sim N(\mathbf{0}, \hat{\sigma}_{u^*}^2 \mathbf{I}_m)$, and $\mathbf{e} \sim N(\mathbf{0}, \hat{\sigma}_e^2 \mathbf{I}_{N^*})$. Then we calculate the corresponding linear predictors $\tilde{\boldsymbol{\eta}}$ and $\log(\tilde{\boldsymbol{y}})$ by using equations

$$\tilde{\boldsymbol{\eta}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\tilde{\boldsymbol{\gamma}} + \mathbf{D}\tilde{\mathbf{u}} \quad \log(\tilde{\boldsymbol{y}}') = \mathbf{X}^* \hat{\boldsymbol{\beta}}^* + \mathbf{Z}^* \tilde{\boldsymbol{\gamma}}^* + \mathbf{D}^* \tilde{\mathbf{u}}^* + \tilde{\mathbf{e}}$$

and $\tilde{\boldsymbol{\pi}}$ and $\tilde{\boldsymbol{y}}'$ are obtained using inverse logit transformation and log-back transformation. The N values of the indicator variable $\tilde{\mathbf{I}}$ are generated performing for each unit a Bernoulli experiment with probability of success equal to the corresponding $\tilde{\pi}_{ij}$. Finally, we obtain the bootstrap population data $(\tilde{y}_{ij}, \mathbf{x}_{ij}, \mathbf{s}_{ij})$ from $\tilde{y}_{ij} = \tilde{I}_{ij} \tilde{y}'_{ij}$.

From this population we drow B bootstrap samples $(\tilde{\boldsymbol{y}}^{(1)}, \dots, \tilde{\boldsymbol{y}}^{(B)})$ maintaining the original small area sample sizes n_i . We refit the model to the bootstrap sample data to obtain the set of bootstrap estimates $\hat{\boldsymbol{\beta}}^{(b)}, \hat{\boldsymbol{\beta}}^{*(b)}, \hat{\sigma}_e^{2(b)}, \hat{\sigma}_u^{2(b)}, \hat{\sigma}_{u^*}^{2(b)}, \hat{\sigma}_\gamma^{2(b)}, \hat{\sigma}_{\gamma^*}^{2(b)}$ and then we estimate the bootstrap means predictors $\widehat{Y}_i^{(b)}$ using (6) and (7).

Indicating with \widetilde{Y}_i the small area mean in the bootstrap population, the estimated root MSEs for the small areas estimates are obtained by:

$$\widehat{\text{RMSE}}_i = \sqrt{\frac{1}{B} \sum_{b=1}^B \left(\widehat{Y}_i^{(b)} - \widetilde{Y}_i\right)^2} \quad (8)$$

3 Application to Real Data

ISTAT drives the Agricultural Census ten-yearly and the sample FSS two-yearly. Both in the Census and in the FSS, the unit of observation is the farm and the data on the amount of land allocated to different crops are collected for each farm. In the FSS, until 2005, the quantity of crops produced was also observed. The FSS is designed to obtain reliable estimates only at regional level, therefore to obtain estimates at sub-regional levels it is necessary to employ SAE estimators using the variables collected at the census time as auxiliary variables. Starting from year 2000 the Agricultural Census georeferences all the farms on the territory, so this information can be used as auxiliary variable as well.

We are interested in producing the mean estimation of grapevine production for the 52 agrarian regions in which Tuscany region is partitioned. The agrarian regions are sub-provincial aggregations of municipalities homogenous respect to natural and agricultural characteristics. The estimates are referred to year 2003 for which the data of the FSS are available. Auxiliary variables and spatial information for each farm referred to 2000 census time. Due to the high correlation values observed over sampled data between the explicative variables at years 2000 and 2003 (about 90 % for the grapevines surface), we suppose that the time lag between the response and the explicative variables should have a negligible effect.

The nature of the study variable does not allow the use of classic small area methods that assume a linear mixed model and don't take into account the spatial structure of the data. A large number of farms don't cultivate grapevines, and a few produce the majority of the total region production. Moreover the cultivation and consequently the production of grapevines for each farm depends on the characteristics of the territory in which the farm is located. Finally, the quantity of grapevine produced by the same allocated surface may change, depending on the soil productivity and on the production choices of the farms (relative to the typology and quality of the produced grapevine).

These practical considerations, confirmed by an explorative analysis of the data, motivate our choice of a two-part model: a logit model for the probability of non-zero grapevine production and a conditional log-transformed model for the mean of non-zero grapevine production. The selection of the covariates among several socio-economic variables (including land use information) available at census time follows the indications obtained from a stepwise regression analysis of the data. For the logit model two auxiliary variables are considered: the surface allocated to grapevines in logarithmic scale and a dummy variable that indicates the selling of grapevine related products, both at year 2000. In the conditional log-transformed model we include the same two variables plus the number of working days done by farm family members in year 2000. Moreover, since both the choice to produce or not produce grapevines ($I_{ij} = 1$ or $I_{ij} = 0$) and the conditioned level of production depend on the characteristics of the farm's location, both models include a smooth function of the Universal Transverse Mercator (UTM) geographical coordinates of each farm's administrative center. Regarding the possibility to include into the

model the specific small area random effect, it results significant only in the log-transformed model. Therefore, recalling (3) and (5), our chosen models are:

$$\eta = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} \quad \log(y') = \mathbf{X}^*\boldsymbol{\beta}^* + \mathbf{Z}^*\boldsymbol{\gamma}^* + \mathbf{D}^*\mathbf{u}^* + \mathbf{e}.$$

The splines knots are selected setting $K = 50$ and using the Clara space filling algorithm of Kaufman and Rousseeuw [11]. The two models are estimated separately, with the logit one fitted through the Penalized Quasi-Likelihood method using all the 2450 farms in the 2003 FSS sample, and the log-transformed one fitted by maximizing the restricted log-likelihood and using only the 961 farms with a strictly positive value of grapevines production.

The estimated models parameters (presented in Table 1) are combined with the census values of the 136817 non sampled farms using (7) to obtain the grapevine production predictions. Finally, expression (6) is applied to obtain the predicted agrarian regions means showed in Table 2 and Fig. 1a.

The estimates in Fig. 1a present an evident geographical pattern, with the higher values in the areas belonging to the provinces of Florence and Siena (the well-known area of Chianti) and the lower values in the north mountainous area of the provinces of Massa Carrara and Lucca. Moreover, the two plots in Fig. 1b show that our results are consistent with both the direct estimate and the expert's estimates.¹ The expert's

Table 1 Estimated parameters with 95 % confidence intervals for the two-part model

Logit model			Log-transformed model		
Parameters ^a	Estimate	Confidence interval	Parameters ^a	Estimate	Confidence interval
Fixed effects			Fixed effects		
Intercept	17.2292	(−23.8177; 58.2760)	Intercept	−0.5709	(−25.0077; 23.8660)
X coordinate	0.0710	(−0.8436; 0.9857)	X coordinate	0.4730	(−0.0130; 0.9591)
Y coordinate	−0.3965	(−1.1956; 0.4026)	Y coordinate	−0.0081	(−0.5179; 0.5018)
Log(grapevine surface 2000)	1.9745	(0.9118; 3.0372)	Log(grapevine surface 2000)	1.2694	(1.2059; 1.3328)
Grapevine products selling	1.0636	(0.0358; 2.0915)	Grapevine products selling	0.6701	(0.5163; 0.8239)
			Family members working days	0.0004	(0.0002; 0.0006)
Random effects			Random effects		
σ_γ	0.2124	(0.0225; 2.0018)	σ_γ^*	0.2394	(0.0795; 0.7206)
σ_ϵ	2.9930	(2.9102; 3.0781)	σ_u^*	0.2189	(0.1243; 0.3854)
			σ_ϵ^*	0.8973	(0.8570; 0.9396)

^aIntercept and coordinates coefficients are not significant but required by model structure

¹Statistics are produced using expert information. Data are provided by local authorities that collect experts evaluations on area and yield of different crops (Source: <http://siqua.istat.it>).

Table 2 Agrarian region level estimates of the mean grapevine production (q) with estimated root mean squared error (RMSE)

Agrarian region	n_i	\widehat{Y}_i	\widehat{RMSE}_i	Agrarian region	n_i	\widehat{Y}_i	\widehat{RMSE}_i
Lunigiana Settentrionale	32	4.67	2.40	Colline tra Era e Fine	26	10.91	1.21
Lunigiana Sud-orientale	29	6.29	4.47	Colline dell'Alto Cecina	29	5.65	2.28
Montagna di Massa	21	4.08	0.68	Colline del Monte Pisano	4	3.45	2.89
Lunigiana Sud-occidentale	20	4.19	2.35	Colline del Medio Cecina	14	14.71	6.85
Garfagnana Occidentale	5	3.23	2.26	Pianura di Pisa	59	4.95	0.92
Garfagnana Centrale	18	3.28	1.68	Casentino	30	13.15	1.48
Garfagnana Orientale	18	3.01	1.82	Alto Tevere	12	5.67	1.44
Val di Lima Lucchese	10	2.76	0.76	Valdarno Superiore	62	37.48	4.39
Montagna della Versilia	4	2.13	0.46	Alta Valle Tiberina	11	10.39	2.96
Pianura della Versilia	108	2.36	1.28	Media Val di Chiana	48	20.52	2.73
Pianura di Lucca	112	7.49	0.67	Colline di Arezzo	114	21.65	2.68
Montagna di Pistoia	185	3.18	1.07	Amiata Orientale	12	6.78	1.11
Val di Nievole	110	3.50	2.15	Alta Val d'Elsa	51	106.59	20.56
Ombone Pistoiese	71	8.16	1.39	Colline del Chianti	34	302.53	75.45
Alto Santerno e Lamone	20	4.44	1.64	Colline di Siena	33	21.97	3.30
Montagna di Vallombrosa	16	42.95	12.98	Val d'Arbia	77	92.41	10.74
Colline del Mugello	40	43.80	7.27	Alta Val di Chiana	63	78.39	18.74
Medio Valdarno	29	51.92	42.50	Val d'Orcia	50	25.77	11.83
Colline di Firenze	34	23.89	3.04	Amiata Occidentale	40	7.15	3.39
Val d'Elsa Inferiore	48	105.60	50.75	Colline dell'Ombone	77	13.61	5.92
Colline del Greve e Pesa	56	225.90	60.02	Colline del Fiora	92	32.90	3.73
Incisa in Val d'Arno	10	45.36	13.22	Colline di Follonica	80	14.57	3.00
Pianura di Fucecchio	27	41.42	3.22	Colline dell'Albenga	117	39.43	4.29
Colline di Livorno	45	8.52	1.43	Pianura di Grosseto	60	19.19	14.43
Colline di Piombino	95	22.35	3.34	Alto Bisenzio	6	3.12	1.70
Valdarno Inferiore	65	20.38	11.66	Colline di Prato	21	13.07	2.89

estimates, produced by ISTAT, are obtained by determination of a crop-specific coefficient of soil productivity and are released at provincial level. We calculated the agrarian region level expert's estimate by multiplying the agrarian region grapevine surfaces at year 2000 with the coefficient of soil productivity at regional level.

Finally, the accuracy of our estimates is assessed by comparing the mean of the area-specific coefficients of variation (CVs) of both the direct and the model based estimates. The CVs are obtained as the ratio between the \widehat{RMSE}_i and the \widehat{Y}_i and this comparison shows an average efficiency gain for our estimates of about 32 %.

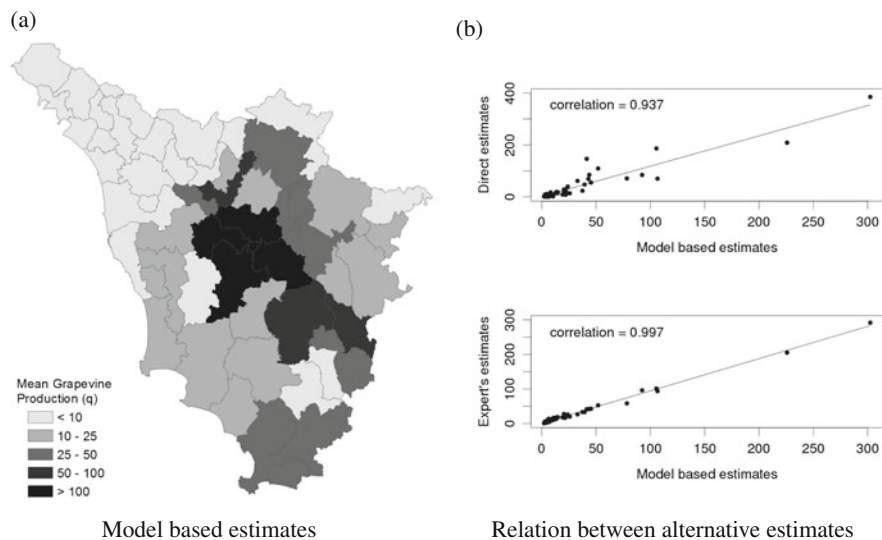


Fig. 1 Agrarian region level estimates of the mean grapevine production. (a) Model based estimates. (b) Relation between alternative estimates

4 Final Remarks

The interest in spatial data analysis is increased in every area of statistical research. Particular interest is given to the possible ways in which spatially referenced data can support local policy makers. Geographical information is frequently available in many areas of observational sciences, and the use of specific techniques of spatial data analysis can improve our understanding of the studied phenomena. Moreover, it is recurrent, not only in agricultural field but also in many other applications such as environmental and economic ones, to encounter variables that have a proportion of values equal to zero and a continuous, often skewed, distribution among the remaining values. The two-part model represents the leading model suggested in literature for this sort of variables. However, there seems to be no studies which combine jointly SAE, models for overdispersed or zero-inflated data and spatial data.

Motivated by the real problem of estimating the per-farm average grapevine production in Tuscany (Italy) at agrarian region level, we have developed a two-part geoadditive model under the framework of SAE. The two-part model provides the flexibility to model data in accordance with a scientifically plausible data generating mechanism and the results are encouraging. Further research should investigate the use of a two-part small area geoadditive model with a correlation between the random effects of the two parts of the model. Such situation leads to a likelihood that does not factor in two separate components, that is the two models cannot be

fitted separately and the estimation process should be approached using a Bayesian approach.

Finally we would like to underline that in literature the application of the two-part model mainly concern biomedical data, however, our results show how this kind of model could be usefully employed in many other application fields.

References

1. Battese, G.E., Harter, R.M., Fuller, W.A.: An error component model for prediction of county crop areas using survey and satellite data. *J. Am. Stat. Assoc.* **83**, 28–36 (1988)
2. Bocci, C., Petrucci, A., Rocco, E.: Geographically weighted regression for small area estimation: an agricultural case study. In: *Proceedings of the XLIII Scientific Meeting of the Italian Statistical Society*, CLEUP, pp. 615–618 (2006)
3. Chandra, H., Chambers, R.: Small area estimation under transformation to linearity. *Surv. Methodol.* **37**, 39–51 (2011)
4. Chandra, H., Chambers, R.: Small area estimation for skewed data in presence of zeros. *Calcutta Stat. Assoc. Bull.* **63**, 249–252 (2011)
5. Chandra, H., Sud, U.: Small area estimation for zero-inflated data. *Commun. Stat. Simul. Comput.* **41**, 632–643 (2012)
6. Cressie, N.: *Statistics for Spatial Data* (revised edition). Wiley, New York (1993)
7. Fay, R.E., Herriot, R.A.: Estimation of income from small places: an application of James-Stein procedures to census data. *J. Am. Stat. Assoc.* **74**, 269–277 (1979)
8. Gosh, P., Albert, P.S.: A Bayesian analysis for longitudinal semicontinuous data with an application to an acupuncture clinical trial. *Comput. Stat. Data Anal.* **53**, 699–706 (2009)
9. Hastie, T.J., Tibshirani, R.: *Generalized Additive Models*. Chapman & Hall, London (1990)
10. Kammann, E.E., Wand, M.P.: Geoadditive models. *Appl. Stat.* **52**, 1–18 (2003)
11. Kaufman, L., Rousseeuw, P.J.: *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, New York (1990)
12. Olsen, M.K., Schafer, J.L.: A two-part random-effects model for semicontinuous longitudinal data. *J. Am. Stat. Assoc.* **96**, 730–745 (2001)
13. Opsomer, J.D., Claeskens, G., Ranalli, M.G., Kauermann, G., Breidt, F.J.: Non-parametric small area estimation using penalized spline regression. *J. R. Stat. Soc. Ser. B* **70**, 265–286 (2008)
14. Pfeiffermann, D., Terryn, B., Moura, F.A.S.: Small area estimation under a two-part random effects model with application to estimation of literacy in developing countries. *Surv. Methodol.* **34**, 235–249 (2008)
15. Ruppert, D., Wand, M.P., Carroll, R.J.: *Semiparametric Regression*. Cambridge University Press, Cambridge (2003)
16. Slud, E.V., Maiti, T.: Small-area estimation based on survey data from a left-censored Fay-Herriot model. *J. Stat. Plann. Inference* **141**, 3520–3535 (2011)
17. Zhang, M., Strawderman, R.L., Cowen, M., Wells, M.E.: Bayesian inference for a two-part hierarchical model: an application to profiling providers in managed health care. *J. Am. Stat. Assoc.* **101**, 934–945 (2006)

A Unified Approach for Defining Optimal Multivariate and Multi-Domains Sampling Designs

Piero Demetrio Falorsi and Paolo Righi

Abstract

The present paper illustrates a sampling method based on balanced sampling, practical and easy to implement, which may represent a general and unified approach for defining the optimal inclusion probabilities and the related domain sampling sizes in many different survey contexts characterized by the need of disseminating survey estimates of prefixed accuracy for a multiplicity both of variables and of domains of interest. The method, depending on how it is parameterized, can define a standard cross-classified or a multi-way stratified design. The sampling algorithm defines an optimal solution—by minimizing either the costs or the sampling sizes—which guarantees: (1) lower sampling errors of the domain estimates than given thresholds and (2) that in each sampling selection the sampling sizes for all the domains of interest are fixed and equal to the planned ones. It is supposed that, at the moment of designing the sample strategy, the domain membership variables are known and available in the sampling frame and that the target variables are unknown but can be predicted with suitable superpopulation models.

1 Introduction

A *unified approach* (UI), which is practical and easy to implement, for defining *optimal multivariate multi-domain sampling* is introduced below.

Some parts of this approach have been described with more details in the papers of Falorsi and Righi [1] and of Righi et al. [2].

P.D. Falorsi • P. Righi (✉)
Istat, Via Balbo 16, Roma, Italy
e-mail: falorsi@istat.it; parighi@istat.it

1. The parameters of interest are $R \times D$ totals, the generic of which, $t_{(dr)} = \sum_{k \in U} y_{rk} \gamma_{dk} = \sum_{k \in U_d} y_{rk}$, represents the total of the variable r ($r = 1, \dots, R$) in the *Domain of Interest* (DI) U_d ($d = 1, \dots, D$) which is a subpopulation (of size N_d) of the population U . The symbols y_{rk} and γ_{dk} denote, respectively, the value of the r -th ($r = 1, \dots, R$) variable of interest of the k -th population unit and the domain membership indicator being $\gamma_{dk} = 1$ if $k \in U_d$ and $\gamma_{dk} = 0$, otherwise. The γ_{dk} values are known, and available in the sampling frame.
2. In addition to the DIs, the other subpopulations relevant in the approach are the *Planned Domains* (PDs), U_h ($h = 1, \dots, H$), which are subpopulations for which the sample designer wants to plan and to fix in advance the sample sizes so as to control the accuracy of the domain estimates. The PDs are in general defined as subpopulations of the DIs. As described below in Sect. 2, the definition of the PDs allows to implement different sampling designs.
3. The random selection of the sample s is implemented with the cube algorithm [3] respecting the following *balancing equations*: $\sum_{k \in s} \delta_k = \sum_{k \in U} \pi_k \delta_k$ in which, with reference to the unit k , π_k is the inclusion probability and $\delta'_k = (\delta_{1k}, \dots, \delta_{Hk})$ is a vector of indicator variables, available in the sampling frame, being $\delta_{hk} = 1$ if $k \in U_h$ and $\delta_{hk} = 0$, otherwise. The above equations guarantee that in each possible sample selection, the realized sample sizes for the planned domains U_h are fixed and equal to the expected ones. Since the PDs are defined as subpopulations of the DIs, also the latter have planned sample sizes.
4. The unknown y_{rk} values are predicted with a simple *working* model, M , $y_{rk} = \tilde{y}_{rk} + u_{rk}$ in which, \tilde{y}_{rk} and u_{rk} ($k = 1, \dots, N$) denote, respectively, the predictions and the random residuals which have the following model expectations:

$$E_M(u_{rk}) = 0 \forall k; E_M(u_{rk}^2) = \sigma_{rk}^2; E_M(u_{rk}, u_{rl}) = 0 \forall k \neq l, \quad (1)$$

further we assume $\sigma_{rk}^2 = \sigma_r^2 v_k^\tau$ where v_k is an auxiliary variable, σ_r^2 and τ are scalar parameters which we assume as known when planning the sampling design. In practice the scalar parameters have to be estimated from pilot or previous survey data.

5. According to Deville and Tillé [4], an approximation of the Measure of the Accuracy (MA) (e.g., the sampling variance or the anticipated variance) of the balanced sampling may be defined as implicit function of the inclusion probabilities and of the squared residual of a generalized linear regression model linking an appropriate transformation of the target variable (which may be known or predicted) to the auxiliary variables involved in the balancing equations. Taking into account the Horvitz Thompson (HT) estimator, $\hat{t}_{(dr)} =$

$\sum_{k \in U} \tilde{y}_{rk} \gamma_{dk} / \pi_k$ of the totals $t_{(dr)}$ and considering the anticipated variance [5] as measure of accuracy, the MA may be expressed by:

$$MA(\hat{t}_{(dr)}) = AV(\hat{t}_{(dr)} | \boldsymbol{\pi}) \cong E_M E_p(\hat{t}_{(dr)} - t_{(dr)} | \boldsymbol{\pi})^2 = f \sum_{k \in U} \left(\frac{1}{\pi_k} - 1 \right) E_M(\eta_{(dr)k}^2), \tag{2}$$

where E_p denotes the expectation over repeated sampling, $\boldsymbol{\pi}$ is the vector of the inclusion probabilities, $\eta_{(dr)k} = (\tilde{y}_{rk} + u_{rk}) \gamma_{dk} - \pi_k g_{(dr)k}$, $f = N / (N - H)$, $g_{(dr)k} = \boldsymbol{\delta}'_k \mathbf{A}^{-1} \sum_{k \in U} \boldsymbol{\delta}_k (\tilde{y}_{rk} + u_{rk}) \gamma_{dk} (1 - \pi_k)$, being $\mathbf{A} = \sum_{k \in U} \boldsymbol{\delta}_k \boldsymbol{\delta}'_k \pi_k (1 - \pi_k)$.

6. The MA may be expressed with a **general expression** based on **stable generic terms** assuming different forms, according to the chosen MA and to the sampling context

$$MA(\hat{t}_{(dr)}) = f \left[\sum_{k \in U} \frac{\omega_{(dr)k}}{\pi_k} - \sum_{k \in U} \left(\varphi_{(dr)k} + \sum_{i=0}^2 \pi_k^i C_{i(dr)k}(\boldsymbol{\pi}) \right) \right] \tag{3}$$

The stable generic terms $\omega_{(dr)k}$ and $\varphi_{(dr)k}$ are fixed quantities (which may be known or predicted) and the $C_{i(dr)k}(\boldsymbol{\pi})$ ($i = 0, 1, 2$) are functions of the vector $\boldsymbol{\pi}$. For instance, the stable generic terms in the case of (2) are given by

$$\begin{aligned} \omega_{drk} &= (\tilde{y}_{rk}^2 + \sigma_{rk}^2) \gamma_{dk}, \quad \varphi_{(dr)k} = \omega_{(dr)k}, \\ C_{0(dr)k}(\boldsymbol{\pi}) &= 2 \boldsymbol{\delta}'_k \mathbf{A}^{-1} [\tilde{y}_{rk} \gamma_{dk} \mathbf{b}_{\tilde{y}(dr)} + \boldsymbol{\delta}_k \sigma_{rk}^2 \gamma_{dk} (1 - \pi_k)], \\ C_{1(dr)k}(\boldsymbol{\pi}) &= - \left[2 \boldsymbol{\delta}'_k \mathbf{A}^{-1} [\tilde{y}_{rk} \gamma_{dk} \mathbf{b}_{\tilde{y}(dr)} + \boldsymbol{\delta}_k \sigma_{rk}^2 \gamma_{dk} (1 - \pi_k)] \right. \\ &\quad \left. + \boldsymbol{\delta}'_k \mathbf{A}^{-1} [\mathbf{b}_{\tilde{y}(dr)} \mathbf{b}'_{\tilde{y}(dr)} + \sum_{j \in U} \boldsymbol{\delta}_j \boldsymbol{\delta}'_j \gamma_{dj} (1 - \pi_j)^2 \sigma_{rj}^2] \mathbf{A}^{-1} \boldsymbol{\delta}_k \right], \\ C_{2(dr)k}(\boldsymbol{\pi}) &= \boldsymbol{\delta}'_k \mathbf{A}^{-1} [\mathbf{b}_{\tilde{y}(dr)} \mathbf{b}'_{\tilde{y}(dr)} + \sum_{j \in U} \boldsymbol{\delta}_j \boldsymbol{\delta}'_j \gamma_{dj} (1 - \pi_j)^2 \sigma_{rj}^2] \mathbf{A}^{-1} \boldsymbol{\delta}_k, \end{aligned}$$

being $\mathbf{b}_{\tilde{y}(dr)} = \sum_{k \in U} \boldsymbol{\delta}_k \tilde{y}_{rk} \gamma_{dk} (1 - \pi_k)$

The expression (3) is suitable for an automated spreadsheet of the algorithm (see below) defining the optimal inclusion probabilities.

7. The *inclusion probabilities* are defined as a solution of the following optimization problem which guarantees lower sampling errors of the domain estimates.

$$\begin{cases} \text{Min} \left(\sum_{k \in U} \pi_k c_k \right) \\ MA(\hat{t}_{(dr)}) \leq \bar{V}_{(dr)} \quad (d = 1, \dots, D; r = 1, \dots, R) \\ 0 < \pi_k \leq 1 \quad (k = 1, \dots, N) \end{cases} \quad (4)$$

where: $MA(\hat{t}_{(dr)})$ is defined according to (3), $\bar{V}_{(dr)}$ is a fixed quantity which defines the threshold of the measure of accuracy of the estimate $\hat{t}_{(dr)}$, and c_k is the cost for collecting information from the unit k . The dominant term in the formula (3) is $\sum_{k \in U} \omega_{(dr)k} / \pi_k$ while the other addenda give a minor contribution. The algorithm for solving the problem (4) consists of three nested calculation loops. The outer loop fixes the values of the functions $C_{i(dr)k}(\boldsymbol{\pi})$. The inner loop defines the π_k^i values which appear as multiplying factor of the functions $C_{i(dr)k}(\boldsymbol{\pi})$ and then the innermost loop is a modified Chromy algorithm [1] which finds the solution to the minimum constrained problem (4) for given values of $\sum_{i=0}^2 \pi_k^i C_{i(dr)k}(\boldsymbol{\pi})$.

2 Some Examples

As a general rule, in order to define the *optimal inclusion probabilities* for a given sampling strategy, the following operations have to be done:

1. Define the DIs and the related PDs.
2. Define the estimator. The HT estimator is considered in the above section; the generalized regression estimator is introduced in Sect. 3.
3. Define the form of the model (1) for predicting the unknown y_{rk} values.
4. Define the form of the MA (e.g., expression 2) and reformulate it according to the general expression (3) which is suitable for the automation of the algorithm for finding optimal inclusion probabilities.

The theory here illustrated is developed for single stage sampling; however, the approach could be easily extended to consider the case of multistage sampling designs.

Some examples are given below in order to demonstrate how the proposed sampling design could represent a way to generalize in a unified framework some well-known sampling designs. In the following the anticipated variance is taken into account as measure of the accuracy. Consider first the univariate and single-domain case in which $R = D = 1$.

Example 1 Optimal Stratified Sampling Assume that the PDs define a single partition of the population U , so that each PD coincides with a stratum, and suppose

that the predicted values of the variable r ($r = R = 1$) of interest are constant in each stratum with uniform stratum variance, e.g., $\tilde{y}_{rk} = \bar{Y}_{rh}$ and $\sigma_{rk}^2 = \sigma_{rh}^2$ (for $k \in U_h$). In this context the UI defines a stratified simple random sampling without replacement (SSRSWOR) design. If the costs c_k are uniform in each planned domain, that is, $c_k = c_h$ for $k \in U_h$, then the stratum sample sizes are computed according to the optimal allocation [6, Sect. 5.5] in which $n_h \approx N_h \sigma_{rh} / \sqrt{c_h}$. If the costs c_k are uniform for all the units in the population, then the well-known Neyman's allocation is realized with $n_h \approx N_h \sigma_{rh}$. Eventually, if the variances are constant over strata, that is, $\sigma_{rh} = \bar{\sigma}_r$, then the proportional allocation is implemented, resulting $n_h \approx N_h$.

Example 2 Optimal PPS Sampling Assume that there is a single planned domain coinciding with the population U and define the stable terms in (3) as $\omega_{drk} = \sigma_{rk}^2$, $\varphi_{(dr)k} = \omega_{(dr)k}$, $C_{i(dr)k}(\boldsymbol{\pi}) = 0$ ($i = 0, 1, 2$). Then, according to the results given in Särndal et al. [7, Chap 12], the UI defines optimal inclusion probabilities proportional to the squared roots of the measures of the heteroscedasticity: $\pi_k \approx \sqrt{x_k}$.

Let us consider now the multivariate multi-domain case and suppose that the sampling estimates have to be calculated for the domains of three domain types T_l ($l = 1, \dots, 3$) each of which defines a partition of the population of U of cardinality D_l being $D = D_1 + D_2 + D_3$. A demonstration of how the sample size of the interest domains may be obtained by different sampling designs is shown below.

Example 3 The standard approach, here denoted as *cross-classified* or *one-way stratified design*, defines the strata by cross-classifying the modalities of the three domain types.

We can obtain the *one-way stratified design* with the UI, by assuring that the U_h coincide with the strata of the one-way stratified design. Then: $H = D_1 \times D_2 \times D_3$. The vectors δ'_k are defined as $(0, \dots, 1, \dots, 0)$ vectors and each U_h can be defined by a specific intersection of the populations of three domains of interest, one for each domain type.

Furthermore if, for every variable r of interest, the predicted values are constant in each stratum with uniform stratum variance, e.g., $\tilde{y}_{rk} = \bar{Y}_{rh}$ and $\sigma_{rk}^2 = \sigma_{rh}^2$ (for $k \in U_h$), then an SSRSWOR design is implemented. After some algebra the (2) becomes

$$AV\left(\hat{t}_{(dr)} \mid \boldsymbol{\pi}\right) = f \sum_{h \in \Gamma_d} \sigma_{rh}^2 \sum_{k \in U_h} (1/\pi_k - 1) = f \sum_{d=1}^D \sum_{h \in \Gamma_d} \sigma_{rh}^2 N_h (N_h/n_h - 1),$$

since the terms \tilde{y}_{rk} disappear and $\pi_k = \pi_h$ (for $k \in U_h$).

Example 4 Consider the previous situation, in which the U_h coincide with the strata of the one-way stratified design and the predicted values are constant in each stratum. If the model variances are proportional to a known value of some auxiliary

variable, e.g., $\sigma_{rk}^2 = \sigma_r^2 v_k$, then a stratified random sampling without replacement with varying inclusion probabilities design is implemented.

Example 5 The PDs U_h are defined combining all the couples of the domains of the domain types; then $H = (D_1 \times D_2) + (D_1 \times D_3) + (D_2 \times D_3)$ and the δ'_k are defined as vectors with three values equal to one, each in correspondence of one of the three above couples, e.g., $(0, \dots, 01, 0 \dots 0, 1, 0, \dots, 0, 1, 0, \dots, 0)$.

Example 6 Some PDs U_h agree with the domains of one population partitions, for instance, T_1 , and the others U_h are defined combining couples of the remaining domain types T_2 and T_3 . Then: $H = D_1 + (D_2 \times D_3)$ and the δ'_k are defined as vectors having two values equal to one, the first in correspondence to the domains of the partition T_1 and the second in correspondence to the couple of the partitions T_2 and T_3 , e.g., $(0, \dots, 01, 0 \dots 0, 1, 0, \dots, 0)$.

Example 7 The PDs U_h agree with the domains of interest; then $H = D_1 + D_2 + D_3$ and the δ'_k are defined as vectors with three ones each in correspondence to one of the three domain types.

The Examples 3 and 4 describe *one-way* (or *standard*) *stratified designs*, while the remaining Examples 5–7 refer to a *multi-way stratified design*. The choice of the sampling design depends on theoretical and operative reasons. From the operative view point the implementation of the one-way stratified design belongs to the current culture of the survey practitioners and it's implementation is uncomplicated, while the *multi-way* design is seldom adopted for defining the sampling strategies of the actual surveys; however these kinds of designs allow to face a lot of empirical contexts in which the traditional approach fails to achieve the target objectives.

3 Remarks on the Regression Estimator

Consider the case in which for producing the sampling estimates, vectors of auxiliary variables are available for all the population units and suppose that the predictions based on this auxiliary information are those given in model (1). In this context, the estimates of interest may be computed with a generalized modified regression estimator, which may be expressed as [8, p. 20]:

$${}_{greg}\hat{t}_{(dr)} = \sum_{k \in U} \tilde{y}_{rk} \gamma_{dk} + \sum_{k \in S} u_{rk} \gamma_{dk} / \pi_k \quad (r = 1, \dots, R; d = 1, \dots, D) \quad (5)$$

An approximation of the anticipated variance of the estimator (4) under balanced sampling is

$$AV \left({}_{greg}\hat{t}_{(dr)} \mid \boldsymbol{\pi} \right) = E_M \left[f \sum_{k \in U} \left(1 / \pi_k - 1 \right) {}_{greg}\eta_{(dr)k}^2 \right],$$

being ${}_{greg}\eta_{(dr)k} = u_{rk} \gamma_{dk} - \pi_k \delta_k \left[\sum_{j \in U} \delta_j \delta'_j \pi_j (1 - \pi_j) \right]^{-1} \sum_{k \in U} \delta_j u_{rj} \gamma_{dj} (1 - \pi_j)$. The expression of the residuals ${}_{greg}\eta_{(dr)k}$ is equivalent to the expression $\eta_{(dr)k}$ given in formula (2), except for the substitution of the terms $(\tilde{y}_{rk} + u_{rk}) \gamma_{dk}$ with $u_{rk} \gamma_{dk}$. The derivation of the expression of stable generic terms of (3) is straightforward.

4 Remarks on Nonresponse

Suppose that, for different causes, it is impossible to collect the survey variables from some sample units. Only to make the things simple, let us further hypothesize that: (1) the phenomenon of nonresponse is substantially different among the PDs U_h ($h = 1, \dots, H$); (2) the response propensities, θ_k , are roughly constant for the units belonging to the subpopulation U_h , that is, $\theta_k \cong \theta_h$ for $k \in U_h$; (3) when planning the sample design, a quite reliable estimate, say $\tilde{\theta}_h$, of the response propensity of the units belonging to U_h may be obtained from the previous surveys. According to the strategy proposed in Särndal and Lundström [9, expression 6.4], the estimator of the totals of interest is calculated with the *calibration estimator*:

$${}_{cal}\hat{t}_{(dr)} = \sum_{k \in s^*} y_{rk} \gamma_{dk} \lambda_k / \hat{\theta}_k \pi_k \quad (r = 1, \dots, R; \quad d = 1, \dots, D), \quad (6)$$

where: s^* is the sample of respondents; $\hat{\theta}_k$ is the sample estimate of the response probability; $\lambda_k = 1 + \left[\sum_{j \in U} \delta_j - \sum_{j \in s^*} (\pi_j \hat{\theta}_j)^{-1} \delta_j \right]' \left[\sum_{j \in s^*} (\pi_j \hat{\theta}_j)^{-1} \delta_j \delta'_j \right]^{-1} \delta_k$. In the context here described, the response probabilities may be estimated by $\hat{\theta}_k = m_h / n_h$ for $k \in s_h^* = s^* \cap U_h$ being m_h the sample size of s_h^* . Let us note that the stratum response probabilities have been introduced with two different symbols, $\tilde{\theta}_h$ and $\hat{\theta}_h$, since the first is an estimate available when planning the sample design, and the latter is estimated from the current survey data. Under the hypothesis that by calibrating in each PD, the nonresponse bias becomes negligible and considering the response phenomenon as a second phase of sampling, then the MA of (6) may be computed by [9, p. 150]:

$$MA({}_{cal}\hat{t}_{(dr)}) = AV \left({}_{cal}\hat{t}_{(dr)} \mid \boldsymbol{\pi} \right) = AV_{sam} + AV_{NR}$$

in which $AV_{sam} = E_m E_p \left(\sum_{k \in s^*} y_{rk} \gamma_{dk} v_k / \pi_k \mid \boldsymbol{\pi} \right)$ is the anticipated variance of the calibrated estimator in the absence of nonresponse and $AV_{NR} = E_m E_p V_q \left(\sum_{k \in s^*} y_{rk} \gamma_{dk} v_k / \hat{\theta}_k \pi_k \mid \boldsymbol{\pi} \right)$ represents the additional part of variability due to the phenomenon of nonresponse, denoting with $V_q(\cdot)$ the variance of (6) over different sets of respondents.

Let $e_{rk} = y_{rk} - \delta'_k \left(\sum_{j \in U} \delta_j \delta'_j \right)^{-1} \sum_{j \in U} \delta_j y_{rj}$ denote the residual with respect to the regression model in which the variables of interest y_{rk} are regressed with

respect to the auxiliary vectors δ'_k and let $e\sigma_{rk}^2$ indicate the model variance of e_{rk} . By adopting, in the phase of planning the sampling design, the reasonable approximations $\left(e_{rk} \gamma_{dk} - \pi_k \delta'_k \mathbf{A}^{-1} \sum_{j \in U} \pi_j \delta'_j \gamma_{dk} e_{rk} (1/\pi_k - 1)\right)^2 \cong e_{rk}^2 \gamma_{dk}$, and $f \cong 1$, the anticipated variance of (6) may be approximated by $AV\left(\widehat{cal}_{(dr)} \mid \boldsymbol{\pi}\right) = \sum_{k \in U} 1/\pi_k \left(e\sigma_{rk}^2 \gamma_{dk} / \tilde{\theta}_k\right) - \sum_{k \in U_e} \sigma_{rk}^2 \gamma_{dk}$, being $\tilde{\theta}_k = \tilde{\theta}_h$ for $k \in U_h$. Thus, having reliable estimates of the response propensities $\tilde{\theta}_k$ and of the model variances $e\sigma_{rk}^2$, it is possible to define the inclusion probabilities that individuate the minimum cost solution, taking into account the additional part of the variance deriving from the expected nonresponse. After some simple algebra, in this context, the terms of the general expression (3) of the MA are given by: $\omega_{(dr)k} = e\sigma_{rk}^2 \gamma_{dk} / \tilde{\theta}_k$, $\varphi_{(dr)k} = e\sigma_{rk}^2 \gamma_{dk}$, $C_{i(dr)k}(\boldsymbol{\pi}) = 0$ (for $i = 0, 1, 2$). Let us note that, if the model variance is constant in each PD h (that is, $e\sigma_{rk}^2 = e\sigma_{rh}^2$ for $k \in U_h$), then $\pi_k = n_h / N_h$ and then the MA may be reformulated according to the sound expression [9, pp. 171–172]

$$AV\left(\widehat{cal}_{(dr)} \mid \boldsymbol{\pi}\right) = \sum_{h=1}^H N_h (N_h / \tilde{m}_h - 1) e\sigma_{rh}^2,$$

being $\tilde{m}_h = \tilde{\theta}_h n_h$ the expected number of respondents in U_h .

References

1. Falorsi, P.D., Righi, P.: A balanced sampling approach for multi-way stratification designs for small area estimation. *Surv. Methodol.* **34**, 223–234 (2008)
2. Righi P., Falorsi P.D.: Optimal allocation algorithm for a multi-way stratification design. In: *Proceedings of the Second ITACOSM Conference*, pp. 49–52. Pisa, 27–29 June 2011
3. Deville, J.-C., Tillé, Y.: Efficient balanced sampling: the Cube method. *Biometrika* **91**, 893–912 (2004)
4. Deville, J.-C., Tillé, Y.: Variance approximation under balanced sampling. *J. Stat. Plan. Inference* **128**, 569–591 (2005)
5. Isaki, C.T., Fuller, W.A.: Survey design under a regression superpopulation model. *J. Am. Stat. Assoc.* **77**, 89–96 (1982)
6. Cochran, W.G.: *Sampling Techniques*. Wiley, New York (1977)
7. Särndal, C.E., Swensson, B., Wretman, J.: *Model Assisted Survey Sampling*. Springer, Berlin (1992)
8. Rao, J.N.K.: *Small Area Estimation*. Wiley, New York (2003)
9. Särndal, C.E., Lundström, S.: *Estimation in Surveys with Nonresponse*. Wiley, New York (2005)

Estimation of Poverty Rate and Quintile Share Ratio for Domains and Small Areas

Risto Lehtonen and Ari Veijanen

Abstract

In the article, we consider the estimation of indicators on poverty and social exclusion for population subgroups or domains and small areas. For at-risk-of-poverty rate, we discuss indirect design-based estimators including model-assisted logistic generalized regression estimators and model calibration estimators. Logistic mixed models are used in these methods. For quintile share ratio, indirect model-based percentile-adjusted predictor methods using linear mixed models are considered. Unit-level auxiliary data are incorporated in the estimation procedures. For quintile share ratio, we present a method called frequency-calibration or n-calibration to be used in cases where aggregate level auxiliary data only are available. Design-based direct estimators that do not use auxiliary data and models are used as reference methods. Design bias and accuracy of estimators are evaluated with design-based simulation experiments using real register data maintained by Statistics Finland and semi-synthetic data generated from the EU-SILC survey.

R. Lehtonen (✉)

Department of Social Research, University of Helsinki, P.O. Box 18, Unioninkatu 35, FI-00014 Helsinki, Finland

e-mail: risto.lehtonen@helsinki.fi

A. Veijanen

Statistics Finland, Helsinki, Finland

e-mail: ari.veijanen@stat.fi

1 Introduction

There is increasing demand in Europe and elsewhere for reliable statistics on poverty and social exclusion produced for regions and other population subgroups or domains (e.g. [4]). Small area estimation of indicators on poverty and social exclusion has been recently investigated in research projects funded by the European Commission under the Framework Programmes for Research and Technological Development (FP). AMELI (Advanced Methodology for European Laeken Indicators) included several specialized sub-projects (work packages) and covered a wide range of topics on poverty, social exclusion and social cohesion [15]. A certain sub-project concentrated on developing small area estimation methods of selected poverty indicators [9]. SAMPLE (Small Area Methods for Poverty and Living Condition Estimates) concentrated on identifying and developing new indicators and models for inequality and poverty with attention to social exclusion and deprivation, as well as to develop and implement models, measures and procedures for small area estimation of the traditional and new indicators and models (<http://www.sample-project.eu/>).

Indicators on poverty and social exclusion investigated in AMELI included at-risk-of-poverty rate, relative median at-risk-of-poverty gap, quintile share ratio (QSR) and the Gini coefficient. In this paper, we discuss small area estimation for poverty rate and QSR using methods introduced in Lehtonen et al. [9] and developed further in Veijanen and Lehtonen [18] and Lehtonen and Veijanen [13]. For poverty rate, we investigate design-based methods including indirect model-assisted logistic generalized regression estimators [6, 7, 10, 12] and semi-direct and semi-indirect model calibration estimators [13]. Logistic mixed models are used as assisting models in these methods. For QSR, indirect model-based methods based on percentile-adjusted predictors [18] using linear mixed models are considered. Unit-level auxiliary data are incorporated in the estimation procedure for these methods. For cases where aggregate level auxiliary data only are available, we present a method called frequency-calibration or n-calibration [18]. Design-based direct estimators that do not use auxiliary data and models are used as reference methods. Design bias and accuracy of estimators are examined with design-based simulation experiments using real register data maintained by Statistics Finland and semi-synthetic data generated from the EU-wide SILC survey (Statistics on Income and Living Conditions; [2]).

The article is organized as follows. Estimation for poverty rate is examined in Sect. 2. Methods for QSR are presented in Sect. 3. Section 4 includes conclusions.

2 Estimation of Poverty Rate for Regions

2.1 Generalized Regression Estimators

The finite population is denoted by $U = \{1, 2, \dots, k, \dots, N\}$, where k refers to the label of population element. A domain $U_d, d = 1, \dots, D$, is a subset of U such as a regional population. The number of units in the domain is denoted by N_d . In sample $s \subset U$, the corresponding subset is defined as $s_d = U_d \cap s$; it has n_d observations. The domains are of unplanned type. Inclusion probabilities are π_k and design weights are $a_k = 1/\pi_k$.

In order to account for possible differences between regions, a mixed model incorporates domain-specific random effects $u_d \sim N(0, \sigma_u^2)$ for domain U_d . For a binary y -variable, a logistic mixed model is of the form

$$E_m(y_k | u_d) = P\{y_k = 1 | u_d; \beta\} = \frac{\exp(\mathbf{x}'_k \beta + u_d)}{1 + \exp(\mathbf{x}'_k \beta + u_d)},$$

where $\mathbf{x}_k = (1, x_{1k}, \dots, x_{pk})'$ is a known vector value for every $k \in U$ and $\beta = (\beta_0, \beta_1, \dots, \beta_p)'$ is a vector of fixed effects common for all domains. The parameters β and σ_u^2 are first estimated from the data, and estimates \hat{u}_d of the random effects u_d for domains are then calculated. Predictions $\hat{y}_k = P\{y_k = 1 | \hat{u}_d; \hat{\beta}\}$ are calculated for $k \in U_d, d = 1, \dots, D$.

The domain total of a study variable y is defined by

$$t_d = \sum_{k \in U_d} y_k, \tag{1}$$

where y_k denotes the value of the study variable for element k . *Horvitz–Thompson (HT) estimator* of domain total (1) is a direct estimator as it only involves observations from the domain of interest:

$$\hat{t}_d = \sum_{k \in s_d} a_k y_k. \tag{2}$$

The estimator is design unbiased but it can have large variance, especially for small domains. HT does not incorporate any auxiliary data.

Generalized regression (GREG) estimators [12, 17] are assisted by a model fitted to the sample. By choosing different models we obtain a family of GREG estimators with same form but different predicted values [6, 7]. Ordinary GREG estimator

$$\hat{t}_{d;GREG} = \sum_{k \in U_d} \hat{y}_k + \sum_{k \in s_d} a_k (y_k - \hat{y}_k) \tag{3}$$

incorporating a linear fixed-effects regression model is often used to estimate domain totals (1) of a continuous study variable. For a binary or polytomous response variable, a logistic model formulation is often chosen. LGREG (logistic GREG; [10]) estimates the frequency f_d of a class C in each domain. A logistic regression model is fitted to indicators $v_k = I\{y_k \in C\}$, $k \in s$, using the design weights. In the MLGREG estimator [7, 11], we use a logistic mixed model involving fitted values $\hat{p}_k = P\{v_k = 1 \mid \hat{\mathbf{u}}_d; \mathbf{x}_k, \hat{\boldsymbol{\beta}}\}$. The random effects are associated with domains U_d or with larger regions U_r . The MLGREG estimator of the class frequency in U_d is

$$\hat{f}_{d;MLGREG} = \sum_{k \in U_d} \hat{p}_k + \sum_{k \in s_d} a_k (v_k - \hat{p}_k). \quad (4)$$

The calculation of \hat{p}_k for all $k \in U_d$, $d = 1, \dots, D$, requires access to unit-level population data on auxiliary variables.

2.2 Model Calibration Estimators

In classical *model-free calibration* [3, 16], a calibration equation is imposed: the weighted sample totals of auxiliary variables reproduce the known population totals. Model-free calibration does not assume an explicit model. In *model calibration* introduced by Wu and Sitter [19], a model is first fitted to the sample. Calibration weights are determined using the fitted values instead of the original auxiliary variables: the weighted sample total of fitted values reproduces the population total of predictions. Our calibration equations for domain estimation specify that the weighted total of fitted values over a subgroup of the sample equals the sum of predictions over the corresponding population subgroup [8, 13].

A model calibration procedure for domain estimation consists of two phases, the *modelling phase* and the *calibration phase*. There is much flexibility in both phases. We have chosen a mixed model formulation involving components that account for spatial heterogeneity in the population. The predictions calculated in the modelling phase are used in the calibration phase when constructing calibration equation and a calibrated domain estimator. In this phase, the target variables for calibration are determined, possibly also including, for example, some of the x -variables used in the modelling phase (e.g. [14]). Calibration can be defined at the population level, at the domain level or at an intermediate level, for example, at a regional level (neighbourhood) that contains the domain of interest. Further, in the construction of the calibrated domain estimator, a “semi-direct” approach involves using observations only from the domain of interest, whereas in a “semi-indirect” approach, also observations outside the domain of interest are included.

In *population-level calibration* [19], the weights must satisfy calibration equation

$$\sum_{i \in s} w_i z_i = \sum_{i \in U} z_i = \left(N, \sum_{i \in U} \hat{y}_i \right)^T, \quad (5)$$

where $z_i = (1, \hat{y}_i)'$. Using the technique of Lagrange multiplier (λ), we minimize

$$\sum_{k \in s} \frac{(w_k - a_k)^2}{a_k} - \lambda' \left(\sum_{i \in s} w_i z_i - \sum_{i \in U} z_i \right)$$

subject to the conditions (5). The equation is minimized by weights

$$w_k(\lambda) = a_k (1 + \lambda' z_k), \quad (6)$$

where $\lambda' = \left(\sum_{i \in U} z_i - \sum_{i \in s} a_i z_i \right)' \left(\sum_{i \in s} a_i z_i z_i' \right)^{-1}$.

In domain estimation, these weights are applied over a domain: the estimator is

$$\hat{f}_{d;pop} = \sum_{k \in s_d} w_k y_k \quad (7)$$

A straightforward generalization of the population-level calibration equation is a domain-level calibration equation

$$\sum_{i \in s_d} w_{di} z_i = \sum_{i \in U_d} z_i = \left(N_d, \sum_{i \in U_d} \hat{y}_i \right)^T, \quad (8)$$

where the weights w_{di} are specific to the domain. From (8) we see that the domain sizes must be known. We minimize

$$\sum_{k \in s_d} \frac{(w_{dk} - a_k)^2}{a_k} - \lambda'_d \left(\sum_{i \in s_d} w_{di} z_i - \sum_{i \in U_d} z_i \right)$$

subject to (8). The solution is $w_{dk} = w_k(\lambda_d)$, defined by (6) for

$$\lambda'_d = \left(\sum_{i \in U_d} z_i - \sum_{i \in s_d} a_i z_i \right)' \left(\sum_{i \in s_d} a_i z_i z_i' \right)^{-1}.$$

The domain estimator is then a weighted domain sum

$$\hat{f}_{d;s} = \sum_{k \in s_d} w_{dk} y_k. \quad (9)$$

We call this estimator *semi-direct*, as the sum only contains y -observations from the domain of interest. It is not a direct estimator, however, as the weights are determined by a model that is fitted to the whole sample. Various semi-direct calibration estimators can be constructed; see Lehtonen and Veijanen [13].

We introduce next various *semi-indirect* estimators. They are weighted sums over a set that is larger than the domain of interest. Our goal is to “borrow strength” from other domains, in an attempt to reduce mean squared error. A semi-indirect domain estimator incorporates whole sample, an enclosing aggregate of regions in a hierarchy of regions or the set of neighbouring domains, including the domain itself. A neighbourhood of a region comprises regions that share a common border with the specified region. In a semi-indirect estimator, we use supersets $C_d \supset U_d$ of domains with corresponding sample subsets $r_d = C_d \cap s$. In our simulations the supersets are composed of domains (at least two neighbour domains are specified for each domain of interest). We define the domain estimator as a weighted sum of all observations in r_d :

$$\hat{f}_{d;r} = \sum_{k \in r_d} w_{dk} y_k \quad (10)$$

The calibration equation is

$$\sum_{i \in r_d} w_{di} z_i = \sum_{i \in U_d} z_i \quad (11)$$

Note that the sum on the left side of (11) extends over r_d which corresponds to population subset C_d , a larger set than U_d on the right side of the equation. We have required that the weights w_{dk} are close to weights a_k in the domain and close to zero outside the domain. The weights minimize

$$\sum_{k \in r_d} \frac{(w_{dk} - I_{dk} a_k)^2}{a_k}$$

where $I_{dk} = I\{k \in s_d\}$, subject to the calibration equations (11) when

$$w_{dk} = I_{dk} a_k + \lambda'_d a_k z_k; \quad \lambda'_d = \left(\sum_{i \in U_d} z_i - \sum_{i \in r_d} I_{di} a_i z_i \right)' \left(\sum_{i \in r_d} a_i z_i z_i' \right)^{-1}.$$

Variance estimation of GREG estimators (3) and (4) can be handled analytically [12] but there is not yet theory of variance estimation of model calibration estimators for domains, so bootstrap is recommended [5].

2.3 Estimation of At-Risk-of-Poverty Rate

At-risk-of-poverty rate is the proportion of poor people in a domain with equivalized income at or below the poverty line t . Our goal is to estimate the domain poverty rate $R_d = (1/N_d) \sum_{k \in U_d} I \{y_k \leq 0.6M\}$. An estimate \hat{M} of reference median income M is obtained from the HT estimated distribution function $\hat{F}_U(t) = (1/\hat{N}) \sum_{k \in s} a_k I \{y_k \leq t\}$. The distribution function defined in domain U_d is estimated by HT:

$$\hat{F}_d(t) = \left(1/\hat{N}_d\right) \sum_{k \in s_d} a_k I \{y_k \leq t\}, \quad \text{where } \hat{N}_d = \sum_{k \in s_d} a_k.$$

Direct (default) *HT-CDF estimator* of poverty rate is

$$\hat{r}_{d;HT} = \hat{F}_d \left(0.6\hat{M}\right). \tag{12}$$

To estimate domain poverty rate by MLGREG or model calibration, we first estimate the domain total of a *poverty indicator* $v_k = I \{y_k \leq 0.6\hat{M}\}$, which equals 1 for persons with income below or at the poverty line and 0 for others. The estimate of the domain total t_d is then divided by the known domain size N_d (or, its estimate \hat{N}_d). For example, the MLGREG estimator of the poverty rate is

$$\hat{r}_{d;MLGREG} = \hat{f}_{d;MLGREG}/N_d. \tag{13}$$

2.4 Simulation Experiments

For design-based simulation experiments, an artificial population of 1 million persons was constructed from real income data of Statistics Finland for seven NUTS level 3 regions in Western Finland. In the simulations, $K = 1000$ samples of $n = 500$ persons were drawn with without-replacement probability proportional to size (PPS) sampling from the unit-level population. For PPS, an artificial size variable was generated as a function of the socio-economic status of household head (wage and salary earners, farmers, other entrepreneurs, pensioners and others). People with low income appear in samples with larger probability than people with large income.

Our models incorporated the following auxiliary variables: age class (0–15, 16–24, 25–49, 50–64, 65 years), gender with interactions with age class, labour force

status (employed, unemployed and not in workforce) and the PPS size variable, in order to account for the unequal probability sampling design used. We created indicators for each class of a qualitative variable. As domains we used the 36 NUTS4 regions. The NUTS classification is hierarchical: each NUTS4 region is contained within a larger NUTS3 region.

The design bias and accuracy of the methods were measured by absolute relative bias (ARB) and relative root mean squared error (RRMSE):

$$ARB = \left| (1/K) \sum_{k=1}^K (\hat{\theta}_{dk} - \theta_d) \right| / \theta_d \quad \text{and} \quad RRMSE = \sqrt{ (1/K) \sum_{k=1}^K (\hat{\theta}_{dk} - \theta_d)^2 } / \theta_d$$

We have chosen GREG with logistic mixed model (MLGREG) and the semi-indirect variant (SI-regional) of model calibration methods for comparison. SI-regional uses the same logistic mixed model as MLGREG. In SI-regional, the superset C_d for a NUTS4 region contains the neighbouring regions and the region itself.

The direct estimator acts as a reference. We present the averages of ARB and RRMSE over domain classes defined by expected domain sample size (Table 1). The logistic mixed model contains regional random effects associated with NUTS4 regions.

Both model-assisted methods indicated larger design bias than the direct estimator, in the smallest domains in particular. In these domains, the semi-indirect model calibration estimator, SI-regional, had larger bias than the LGREG estimator MLGREG. The large bias probably was caused by the possible heterogeneity of the ensemble of domains forming a superset C_d . The bias declined rapidly with increasing domain sample size. With respect to accuracy, both model-assisted methods clearly outperformed the direct estimator. SI-regional showed better accuracy than MLGREG, but the difference vanished when domain sample size increased. The choice of the model did not have much effect on the model-assisted methods.

Table 1 Mean absolute relative bias (ARB) (%) and mean relative root mean squared error (RRMSE) (%) of three poverty rate estimators over domain size classes, under the PPS sampling design with sample size $n = 500$

Estimator	Expected domain sample size				All
	<5	5–9	10–49	>49	
<i>Mean ARB (%)</i>					
Direct	6.0	2.7	2.1	1.2	2.9
MLGREG	7.7	2.9	1.7	1.8	3.1
SI-regional	16.6	9.4	2.2	1.8	7.8
<i>Mean RRMSE (%)</i>					
Direct	123.7	94.5	64.4	29.8	85.8
MLGREG	119.5	88.3	62.2	30.1	81.4
SI-regional	99.7	80.7	59.1	30.0	73.9

GREG and MC are assisted by logistic mixed model with region-level random effects

3 Estimation of Quintile Share Ratio for Regions

In an indirect model-based predictor-type estimator for QSR based on unit-level auxiliary data, predictions obtained from a linear mixed model are plugged into the formula of QSR defined at the population level. The estimator is expected to have small variance but as a model-based estimator, it can be seriously design biased. To decrease design bias, Veijanen and Lehtonen [18] defined a transformation that reduces the differences between the percentiles of transformed predictions and the percentiles of sample values.

QSR is the ratio of the average equivalized income in the poorest quintile to the average in the richest quintile. Each quintile contains 20 % of people in the population. The QSR is usually calculated as a ratio of Hájek estimators, weighted means over quintiles comprising 20 % of design weights. It is a *direct* estimator. The predictor-type estimator of quintile share in a domain is the ratio of averages of predictions in the first and fifth quintiles.

3.1 Percentile-Adjusted Predictions

Differences between domains are described by incorporating domain-specific random effects $u_d \sim N(0, \sigma_u^2)$ in a linear mixed model $Y_k = \mathbf{x}'_k \boldsymbol{\beta} + u_d + \varepsilon_k, k \in U_d, \varepsilon_k \sim N(0, \sigma^2)$. The parameters $\boldsymbol{\beta}, \sigma_u^2$ and σ^2 are estimated from the data and the values of the random effects are predicted. The predictions are calculated as $\hat{y}_k = \mathbf{x}'_k \hat{\boldsymbol{\beta}} + \hat{u}_d, k \in U_d$.

As the distribution of income is right-skewed, a linear mixed model is fitted to $z_k = \log(y_k + 1)$ and the predictions \hat{z}_k are back-transformed to $\hat{y}_k = \exp(\hat{z}_k) - 1$. After this transformation, the distribution of predictions is usually still far more concentrated around the average than the distribution of observed values $y_k (k \in s_d)$. This leads to too small estimates of QSR. We measure the difference in distribution by percentiles, denoted \hat{p}_{cd} for predictions and p_{cd} for sample values. The percentiles p_{cd} are obtained from the estimated c.d.f. $\hat{F}_d(y) = (1/\hat{N}_d) \sum_{k \in s_d} a_k I\{y_k \leq y\}$. We obtain new predictions with better distribution by applying a transformation $\tilde{y}_k = e^{\alpha_d \hat{y}_k^\gamma}$ with parameters α_d and γ that are chosen to minimize the sum of differences of logarithms of percentiles,

$$\sum_d \sum_{c=1}^C (\log(\tilde{p}_{cd}) - \log(p_{cd}))^2 = \sum_d \sum_{c=1}^C (\alpha_d + \gamma \log(\hat{p}_{cd}) - \log(p_{cd}))^2,$$

where \tilde{p}_{cd} denote the percentiles of \tilde{y}_k and $C = 99$. The OLS estimates of α_d and γ are incorporated in *percentile-adjusted*, or *p-adjusted*, predictions defined by

$$\log(\tilde{y}_k) = \hat{\alpha}_d + \hat{\gamma} \log(\hat{y}_k) \quad (k \in U_d). \tag{14}$$

In the experiments, the transformation (14) was applied using percentiles \hat{p}_{cd} and p_{cd} calculated only from positive values of predictions and sample observations. In some recent simulation experiments, we have preferred using n_d (domain sample size) equally spaced quantiles in each domain instead of calculating the percentiles for every $c = 1, 2, \dots, 99$. For MSE estimation, different variants of bootstrap can be used.

3.2 Frequency-Calibrated Predictor

Unit-level information about the population is not always available. Suppose all the auxiliary variables are qualitative and we know merely their marginal class frequencies in each domain. Then it is seemingly impossible to calculate a predictor incorporating nonlinear predictions, and the only alternatives are typically linear estimators, such as GREG, and area-level estimators. The domain sums of nonlinear predictions are tractable only if the joint frequencies of auxiliary variables are known in each domain. But could we estimate the unknown joint frequencies? Veijanen and Lehtonen [18] introduce a working solution under the constraint of given marginal frequencies.

The joint domain frequencies of auxiliary variables are equivalent with the domain frequencies of distinct values of the vectors \mathbf{x}_k . Their components include the constant and the class indicators. Let us denote the set of distinct values of \mathbf{x}_k , $k \in s_d$, by $X_d = \{z_1, z_2, \dots, z_m\}$. The unknown population frequency n_z of $z \in X_d$ is first estimated by HT as $\hat{n}_z = \sum_{k \in s_d} a_k I\{\mathbf{x}_k = z\}$. We might base our predictor on these estimates, but we want to use information about the known marginal frequencies over population domains, written as equation

$$\sum_{k \in U_d} \mathbf{x}_k = \sum_{z \in X_d} n_z z = t_d, \quad (15)$$

where t_d contains the marginal frequencies, that is, the domain sums of corresponding class indicators, and the domain size corresponding to the constant included in \mathbf{x}_k . By calibration, we derive frequencies \tilde{n}_z that satisfy the calibration equations (15) without deviating too far from the HT estimates \hat{n}_z . The distance of $\tilde{n} = (\tilde{n}_z; z \in X_d)$ to $\hat{n} = (\hat{n}_z; z \in X_d)$, $\sum_{z \in X_d} (1/\hat{n}_z) (\hat{n}_z - \tilde{n}_z)^2$, is minimized subject to the calibration equation (15) by

$$\tilde{n}_z = \hat{n}_z (1 + \lambda'_d z), \quad (16)$$

where the Lagrange multiplier λ_d is defined by

$$\lambda'_d = \left(t_d - \sum_{z \in X_d} \hat{n}_z z \right)' \left(\sum_{z \in X_d} \hat{n}_z z z' \right)^{-1}.$$

The estimated frequencies are used to construct the vector of predictions: for each $z \in X_d$, the vector contains \tilde{n}_z copies of the corresponding fitted value (the frequency estimates are rounded). The predictions are then transformed by (14) and incorporated in the *frequency-calibrated*, or *n-calibrated*, predictor.

There are some practical problems in applying (16). The matrix inversion in the Lagrange multiplier is not possible unless we exclude the indicator variables of classes that do not appear in a sample domain. In addition, if two indicator variables had identical values in the domain, the latter variable was removed. We replaced negative estimates \tilde{n}_z with zeroes, at the expense of violating the calibration equations. This was necessary in 10 % of the estimates. In the experiments, we used the HT estimates \hat{n}_z in the case of other numerical failures.

3.3 Simulation Experiments

In simulation experiments, we used a semi-synthetic data set of about 10 million persons, constructed from SILC data sets [1]. It describes realistically the regional and demographic variation of income in the EU. $K = 1000$ samples were drawn by SRSWOR. We made two experiments, one with $n = 200$, the other with $n = 2000$. There were 40 regions, classified by expected sample size in experiments with $n = 200$ to minor (less than 5 units) and major (5–10 units), in experiments with $n = 2000$ to minor (less than 45 units), medium (45–55 units) and major regions (more than 55 units). We modelled the equivalized income using age class and gender with interactions, attained education level (ISCED), activity (working, unemployed, retired or otherwise inactive) and degree of urbanization of residence (three classes). We applied ML in fitting the mixed models with random effects associated with regions. Design bias and accuracy were described by ARB and RRMSE (see Sect. 2).

The p-adjusted predictor based on (14) was much more accurate than the direct estimator in all domain size classes (Table 2). However, the estimator was design biased, especially in small domains. In domains with less than 10 units, the p-adjusted predictor yielded estimates on an average about 40 % smaller than the true values. On the other hand, the direct estimator was even more biased in the smallest domains. As expected, the frequency-calibrated estimator (Eq. 16) was less accurate than the p-adjusted predictor but more accurate than the direct estimator. In the smallest domains, however, numerical problems with calibration were too common, and the estimates were deemed unusable.

Table 2 Mean absolute relative bias (ARB) (%) and mean relative root mean squared error (RRMSE) (%) of quintile share ratio (QSR) estimators over domain size classes in two different sample settings, under linear mixed model formulation with region-level random effects

Estimator	Sample size $n = 200$		Sample size $n = 2000$		
	Expected domain sample size		Expected domain sample size		
	<5	5–10	<45	45–55	>55
<i>Mean ARB (%)</i>					
Direct	125.0	91.5	4.9	4.6	3.4
p-adjusted predictor	43.9	40.6	12.3	8.6	5.7
n-calibrated predictor	(not used)	(not used)	11.1	13.3	10.6
<i>Mean RRMSE (%)</i>					
Direct	250.3	212.5	43.5	41.7	38.5
p-adjusted predictor	49.7	46.4	16.0	13.6	11.4
n-calibrated predictor	(not used)	(not used)	31.3	29.6	25.9

4 Conclusions

In the estimation of at-risk-of-poverty rate for regions, the GREG estimator and the semi-indirect model calibration estimator outperformed the direct estimator with respect to bias and accuracy. Of the model-assisted estimators based on the same underlying logistic mixed model, the model calibration estimator resulted in better accuracy but slightly larger bias than the GREG estimator, in the smallest domains in particular. The difference vanished when domain sample size increased.

For QSR, we presented certain model-based estimators as alternatives to the direct estimator. The method based on percentile-adjusted predictions decreased the design bias and improved the accuracy, when compared with the direct estimator. Unit-level auxiliary information was used in these methods. The n-calibrated predictor, which used aggregate level auxiliary information, also outperformed the direct estimator.

References

1. Alfons, A., Templ, M., Filzmoser, P., Kraft, S., Hulliger, B., Kolb, J.-P., Münnich, R.: Report on outcome of the simulation study. Research Project Report WP6 (D6.2, FP7-SSH-2007-217322 AMELI). Available at: <http://svn.uni-trier.de/AMELI> (2011)
2. Atkinson, A.B., Marlier, E. (eds.): Income and Living Conditions in Europe. European Commission, Publications Office of the European Union, Luxembourg (2010)
3. Deville, J.-C., Särndal, C.-E.: Calibration estimators in survey sampling. *J. Am. Stat. Assoc.* **87**, 376–382 (1992)
4. European Commission: Combating poverty and social exclusion. A statistical portrait of the European Union 2010. Publications Office of the European Union, Luxembourg (2010)

5. Gershunskaya, J., Jiang, J., Lahiri, P.: Resampling methods in surveys. In: Rao, C.R., Pfeffermann, D. (eds.) *Handbook of Statistics. Sample Surveys. Inference and Analysis*, vol. 29B, pp. 219–249. Elsevier, Amsterdam (2009). Chapter 28
6. Lehtonen, R., Särndal, C.-E., Veijanen, A.: The effect of model choice in estimation for domains, including small domains. *Surv. Methodol. J.* **29**, 33–44 (2003)
7. Lehtonen, R., Särndal, C.-E., Veijanen, A.: Does the model matter? Comparing model-assisted and model-dependent estimators of class frequencies for domains. *Stat. Transit.* **7**, 649–673 (2005)
8. Lehtonen, R., Särndal, C.-E., Veijanen, A.: Model calibration and generalized regression estimation for domains and small areas. In: *SAE 2009 Conference on Small Area Estimation*, Elche, Spain, June–July 2009. Available at: <http://icio.umh.es/congresos/sae2009/> (2009)
9. Lehtonen, R., Veijanen, A., Myrskylä, M., Valaste, M.: Small area estimation of indicators on poverty and social exclusion. Research Project Report WP2 (D2.2, FP7-SSH-2007-217322 AMELI). Available at: <http://svn.uni-trier.de/AMELI> (2011)
10. Lehtonen, R., Veijanen, A.: Logistic generalized regression estimators. *Surv. Methodol. J.* **24**, 51–55 (1998)
11. Lehtonen, R., Veijanen, A.: Domain estimation with logistic generalized regression and related estimators. In: *IASS Satellite Conference on Small Area Estimation*, Riga: Latvian Council of Science, 121–128 (1999)
12. Lehtonen, R., Veijanen, A.: Design-based methods of estimation for domains and small areas. In: Rao, C.R., Pfeffermann, D. (eds.) *Handbook of Statistics. Sample Surveys. Inference and Analysis*, vol. 29B, pp. 219–249. Elsevier, Amsterdam (2009). Chapter 31
13. Lehtonen, R., Veijanen, A.: Small area poverty estimation by model calibration. *J. Indian Soc. Agric. Stat.* **66**, 125–133 (2012)
14. Montanari, G.E., Ranalli, M.G.: Multiple and ridge model calibration. In: *Proceedings of Workshop on Calibration and Estimation in Surveys 2009*. Statistics Canada (2009)
15. Münnich, R., Zins, S., Alfons, A., Bruch, C., Filzmoser, P., Graf, M., Hulliger, B., Kolb, J.-P., Lehtonen, R., Lussmann, D., Meraner, A., Myrskylä, M., Nedyalkova, D., Schoch, T., Templ, M., Valaste, M., Veijanen, A.: Policy Recommendations and Methodological Report. Research Project Report WP10 (D10.1/D10.2, FP7-SSH-2007-217322 AMELI). Available at: <http://svn.uni-trier.de/AMELI> (2011)
16. Särndal, C.-E.: The calibration approach in survey theory and practice. *Surv. Methodol.* **33**, 99–119 (2007)
17. Särndal, C.-E., Swensson, B., Wretman, J.: *Model Assisted Survey Sampling*. Springer, New York (1992)
18. Veijanen, A., Lehtonen, R.: Percentile-adjusted estimation of poverty indicators for domains under outlier contamination. *Stat. Transit.* **12**, 345–356 (2011)
19. Wu, C., Sitter, R.R.: A model-calibration approach to using complete auxiliary information from survey data. *J. Am. Stat. Assoc.* **96**, 185–193 (2001)

A Sample Survey on Inactive Students: Weighting Issues in Modelling the Inactivity Status

Lucio Masserini and Monica Pratesi

Abstract

In this paper we analyze the issue of inactive university students, that is, students with zero university credits in career for at least a calendar year. A focus on this topic is important not only for the negative effects of inactivity on students and their families but also because an increasingly amount of the Ordinary Financing Fund (FFO) is allocated to the universities taking into account the performance of the students' career. Data were collected through a CATI questionnaire administered to a stratified simple random sample of 1945 students enrolled at the University of Pisa in the academic year 2010–2011. The probability of being in the inactive status is modelled specifying a two-level random intercepts logistic regression model, after dealing with the issue of weighing the statistical units.

1 Introduction

The presence and the reasons of a number of inactive students in the Italian universities are not widely investigated topics. These students spend at least one year without achieving university credits and represent a direct cost for themselves and their families besides creating further burdens for the university system. In fact, the reducing Ordinary Financing Fund (FFO) and the increasingly portion allocated taking into account the quality of the educational offer and the results of the

L. Masserini (✉)

Statistical Observatory of the University of Pisa, Lungarno Pacinotti 43, Pisa 56126, Italy
e-mail: l.masserini@adm.unipi.it

M. Pratesi

Department of Statistics and Applied Mathematics to Economics, University of Pisa, Pisa, Italy
e-mail: m.pratesi@ec.unipi.it

training processes require a focus on this topic. More specifically, the Ministry for Education, University and Research (MIUR) adopted two indicators for allocating a portion of the FFO rewarding share in 2012. The first (A1) refers to regular students in the academic year 2010–2011 who achieved at least five university credits during the calendar year 2011; the second (A2) concerns the ratio between the credits achieved in the calendar year 2011 and the theoretical ones by students enrolled in the academic year 2010–2011. From this point of view, it is relevant to characterize students having periods of inactivity during their university career. Frequently, a number of these students spend more than one year being in an inactive status and it cannot be excluded that such a period may result in a later dropout. The literature on inactive university students is rather poor. The only research known to the authors is carried out by the Faculty of Philosophy of the University of Rome “La Sapienza” and is limited to the students enrolled in the degree programme “Education and Training” [1]. Other studies focused on the related issue of dropouts [2]. A brief descriptive analysis, referred to the academic year 2008–2009, can be found in the “Eleventh Annual Report on the Condition of the University System” worked out by the National Committee for the Evaluation of the University System [3].

In this paper we aim at modelling the probability of being in the inactive status in the calendar year 2011 for students enrolled at the University of Pisa in the academic year 2010–2011. We describe the sample survey designed by the Statistical Observatory of the University of Pisa and then we specify a two-level random intercepts logistic regression model for the probability of being in the inactive status (Sect. 2). This section describes the data collection, the statistical model, and the variables. Section 2.1 focuses on the sampling design and the questionnaire. Section 2.2 specifies a two-level random intercepts logistic regression model and discusses the estimation choice between model-based and design-based approach. Section 2.3 describes the variables. Finally, we present the results (Sect. 3) and the final remarks (Sect. 4).

2 Data and Method

2.1 Data Collection

The data collection process was conceived with reference to the condition of inactive status which defines students with zero university credits (or examinations passed) in a calendar year. A stratified simple random sample of 1945 students was selected from the target population (51,758 enrolled students in the academic year 2010–2011). The stratification criteria were: *Activity status* (Active, Inactive), *Regularity of the enrollment condition* (Regular, Not Regular by 1–2 years, Not Regular by more than 2 years), *Subject area of the course of study* (four areas under the current regulation—Medicine and Health, Science and Mathematics, Social Sciences, and Humanities—and a miscellanea under the old regulation), and *Status of freshmen* (Yes, Not). An approximately equal number of active (968) and inactive (977) students were obtained by selecting the sample units with unequal probabilities

from each of the two different strata. The allocation of the students into the other strata was proportional to the population size. Data were collected using a CATI (Computer Assisted Telephone Interviewing) system. Interviews were carried out by a group of qualified part-time students. The questionnaire was divided into five main sections (a–e): the enrollment condition and the short time perspectives (a); school experience before enrolling at the university (b); motivations for enrollment and choice of the actual course of study (c); experience in the university system (attendance, network with other students, Erasmus programme, stages, tutorships) (d); evaluation of personal dimensions (interest in subjects of study, skills and abilities in the study, etc.) (e); socio-demographic data (f). The survey was conducted from March 20 to May 5, 2012. The average time for interview was about 14 min (standard deviation 3 min), the average number of call attempts to complete an interview was 6.3 (maximum number of call attempts was 15) and the refusal rate was 3.2 %.

2.2 Statistical Model

Data show a typical hierarchical structure [4, 5], in which lower-level units (individuals) are nested within higher-level units (clusters). Here, students are clustered into courses of study that, in turn, are nested within faculties, defining a natural three-level hierarchical data structure. As a consequence, students within the same cluster for each level of this data structure, sharing unobserved factors, usually show correlated values of the response variable. Ignoring this structure and analyzing lower-level units as if they were independent produce biased standard errors of the regression coefficients [6]. The outcome variable is a binary response and distinguishes the students with university credits from the students with zero credits in the calendar year 2011. The analysis is performed using a two-level random intercepts logistic regression model [7, 8] by taking into account the hierarchical data structure and the nature of the response variable. A three-level model was firstly estimated in order to investigate the so-called faculty-effect but the results were not significant and show that only the lowest level of clustering (course of study) affects the responses. Since the characteristics of degree programmes reflect both difficulties due to subjects of study and to teaching organization, they affect the probability of being in the inactive status more than difficulties attributable to the faculties. The latter are usually related to a more general and less defined “environmental” effect, often restricted only to location and logistics of the classrooms and to the other premises where educational activities take place. As a result, the effects due to the hierarchical data structure are limited to those exercised by the courses of study that represent the second-level units.

The model can be defined from a binary outcome $y_{ij} \in (0, 1)$ observed on student i , with $i = 1, 2, \dots, N_j$, in course of study j , with $j = 1, 2, \dots, G$, which takes on the value of 0 for students with credits in 2011 and 1 otherwise. From this outcome, the probability that y_{ij} takes on the value of 1 can be defined as

$P_{ij} = Pr(y_{ij} = 1)$. A typical transformation of this probability is the logit transformation, where the logit is defined in terms of the natural logarithm of the odds ratio, indicated as $\ln(P_{ij}/1 - P_{ij})$. The two-level random intercepts logistic regression model can be expressed using the single equation mixed model formulation [8]:

$$\ln\left(\frac{P_{ij}}{1 - P_{ij}}\right) = \beta_0 + \beta_1 \mathbf{x}_{ij} + \boldsymbol{\gamma} \mathbf{w}_j + u_j,$$

where \mathbf{x}_{ij} is a vector of covariates for student i in course study j , \mathbf{w}_j is a vector of covariates characterizing the course of study j whereas $u_j \sim N(0, \sigma_u)$ is an unobservable quantity shared by the students within a particular course of study capturing all relevant factors not accounted for by the observed covariates. The parameters to be estimated are β_0 , representing the population average of the transformed probabilities, β_1 and $\boldsymbol{\gamma}$, the regression coefficients, and σ_u , the standard deviation of the random effects, whose magnitude indicates the strength of the influence of the specific course of study j .

A challenging issue concerns model estimation because data collection is conceived using a stratified simple random sample, with sampling proportions differing between strata. In presence of a complex survey design, two different approaches to analytic inference are design-based and model-based [9]. In a model-based approach, parameters are biased only if the distributions of the residuals are affected by the sampling design. In this case the survey design is said to be *informative* [10] and the model holding for the sample data is different from the model holding in the population. Failure to account for the effects of informative sampling may yield large biases and erroneous conclusions [11]. In presence of an informative sampling design, parameters estimation can be performed using sampling weights. On this basis, model-based (unweighted) and design-based (weighted) estimates were compared. The weighting methods for hierarchical models are developed by Pfeffermann et al. [12], Grilli and Pratesi [13], Asparouhov [14], and Rabe-Hesketh and Skrondal [15]. Here, design-based estimates are obtained using the full pseudo-maximum-likelihood estimation method proposed by Rabe-Hesketh and Skrondal [15] for generalized linear mixed model and computed by Stata programme *gllamm* using sampling weights. In our sampling design, weights are defined only at student-level because all the courses of study in the population are included in the sample. Weights are rescaled using “Method 2” in Pfeffermann et al. [12]:

$$w_i^* = w_i \frac{n_j}{\sum_j w_{i|j}}$$

Here, $w_i = 1/\pi_i$, and π_i defines the probability that student i is included in the sample. Moreover, n_j is the number of student in course of study j whereas $w_{i|j}$ is the weight for student i in course of study j . For the comparison of each parameter estimates, under model-based and design-based approach, Asparouhov’s

[14] informativeness measure (I_2) is employed:

$$I_2 = \frac{\hat{\theta}_h^M - \hat{\theta}_h^W}{\sqrt{\hat{\sigma}_{hh}^M}},$$

where $\hat{\theta}_h^M$ and $\hat{\theta}_h^W$ are the model-based and the design-based estimator for parameter h , whereas $\hat{\sigma}_{hh}^M$ is the model-based standard error. Values of this measure larger than 2 can be considered unacceptable.

2.3 Variables

Covariates characterizing students and their family status and covariates describing the actual (university) and the previous (school) student experience are: *Age*, *Gender* (0 = Male; 1 = Female), *Having a job* (0 = Not; 1 = Yes), *Marital status* (0 = Unmarried; 1 = Married, Cohabitant, or Divorced), *Having children* (0 = Not; 1 = Yes), *Father and Mother educational qualification* (0 = Undergraduate; 1 = Graduate), *School final mark* (0 = >90/100; 1 = <91/100), *Grammar school* (0 = Not; 1 = Yes), *Enrollment condition* (0 = Regular or Not Regular by 1–2 years; 1 = Not Regular by more than 2 years), *Attending lessons* (0 = Little, Enough, or A lot; 1 = None), *Current course is the one in which you wish to graduate* (0 = Yes; 1 = Not), *Having contacts with classmates outside the university* (0 = Little, Enough, or A lot; 1 = None), *Enrolled to the same course of registration* (0 = Yes; 1 = Not). At course-level, covariates are obtained from the questionnaire for the evaluation of the teaching activities and can be considered as proxies of the course quality: % of *Overall satisfaction for teaching activities*, *Personal tolerability of the study load*, and % of *Overall satisfaction organization of the teaching activities*.

3 Results

Preliminary descriptive statistics shows that at the University of Pisa the population percentage of inactive students in 2011 compared to the students enrolled in the academic year 2010–2011 is 18.4 (9506 out of 51,758). Limited to students enrolled to courses degree under the new regulation this percentage is 17.3 for the students of the first cycle degree but drops to 11.6 among the freshmen. Instead, for the second cycle degree, the percentage is even lower, 9.8.

Data analysis is limited to students under the new regulation and it is carried out by estimating separate models for the first and the second cycle degree courses. We decided to estimate different models because of the considerable differences between undergraduate and postgraduate programmes in terms both of students (for motivations, expectations, and so on) and for subjects of study. Such a difference, as

reflected by the observed inactivity rates, may result not only in different regression coefficients but also for the possible existence and/or magnitude of a course-effect. Students enrolled under the old regulation are excluded from the analysis because their university experience is completely different from that of the other students. Moreover, the small number of observations does not allow for the estimation of a separate model. The results compare the model-based approach with the design-based approach obtained by the full pseudo-maximum-likelihood estimation. Model selection is pursued firstly by estimating the null model in order to quantify the unobserved heterogeneity induced by the clustering of students into degree courses. Then, we proceed by introducing the student-level covariates and finally the model is augmented with the course-level covariates as far as they help in explaining the residual second-level variability. The results, reporting regression coefficients (Coef.), standard error (SE), p -value (P), the transformed probability (Prob), and the Asparouhov's [14] informativeness measure (I_2), are shown, respectively, in Sect. 3.1 for the first cycle degree courses (Tables 1 and 2) and Sect. 3.2 for the second cycle degree courses (Table 3).

Table 1 Parameters estimates for unweighted estimation

Model parameters	Coef.	SE	P	Prob	I_2
<i>Fixed effects</i>					
Constant (reference profile)	-1.318	0.158	0.000	0.268	3.203
Enrolled to the same course of registration: Not vs Yes	0.668	0.151	0.000	0.343	0.178
Repeat years during high school: Yes vs Not	0.514	0.196	0.009	0.309	0.047
Current course is the one in which you wish to graduate: Yes vs Not	-0.629	0.256	0.014	0.125	-0.418
Attends lessons more or less regularly: Not vs Yes	0.791	0.202	0.000	0.371	0.187
Having contacts with classmates outside the university: Not vs Yes	0.519	0.190	0.006	0.310	0.318
Having a job: Yes vs Not	1.170	0.143	0.000	0.463	0.051
Having children: Yes vs Not	0.376	0.141	0.008	0.281	0.061
School final mark $\geq 90/100$	-0.601	0.162	0.000	0.128	-0.427
Grammar school: Yes vs Not	-0.325	0.138	0.018	0.162	0.099
Regular or not regular registration by 1-2 years: Not vs Yes	1.585	0.167	0.000	0.566	0.346
Random effect (course of study standard deviation)	0.332	0.111	0.022	-	-
Courses with high probability of inactivity (-2 standard deviation)	-	-	-	0.121	-
Courses with low probability of inactivity (+2 standard deviation)	-	-	-	0.342	-

Table 2 Parameters estimates for the full pseudo-maximum-likelihood estimation

Model parameters	Coef.	SE	P	Prob
<i>Fixed effects</i>				
Constant (reference profile)	-2.590	0.207	0.000	0.070
Enrolled to the same course of registration: Not vs Yes	0.598	0.164	0.000	0.120
Repeat years during high school: Yes vs Not	0.493	0.190	0.009	0.109
Current course is the one in which you wish to graduate: Yes vs Not	-0.417	0.260	0.109	0.047
Attends lessons more or less regularly: Not vs Yes	0.708	0.240	0.003	0.132
Having contacts with classmates outside the university: Not vs Yes	0.380	0.194	0.047	0.099
Having a job: Yes vs Not	1.151	0.161	0.000	0.192
Having children: Yes vs Not	0.354	0.185	0.056	0.097
School final mark $\geq 90/100$	-0.429	0.216	0.023	0.047
Grammar school: Yes vs Not	-0.362	0.168	0.031	0.050
Regular or not regular registration by 1–2 years: Not vs Yes	1.443	0.198	0.000	0.241
Random effect (course of study standard deviation)	0.257	0.098	–	–
Courses with high probability of inactivity (–2 standard deviation)	–	–	–	0.028
Courses with low probability of inactivity (+2 standard deviation)	–	–	–	0.107

Table 3 Parameters estimates for the full pseudo-maximum-likelihood estimation

Model parameters	Coef.	SE	P	Prob
Constant (reference profile)	-3.075	0.230	0.000	0.044
Attends lessons more or less regularly: Not vs Yes	1.036	0.274	0.002	0.115
Having a job: Yes vs Not	0.780	0.274	0.004	0.091
Having children: Yes vs Not	0.594	0.393	0.031	0.077
Regular or not regular registration by 1–2 years: Not vs Yes	1.897	0.283	0.000	0.235
Registration motivated mainly for job opportunities Yes vs Not	-0.549	0.283	0.052	0.025

3.1 First Cycle Courses Degree

The analysis comprises 1382 students within 139 courses of study. The average number of students for each course is 9.9 with a minimum of 5 and a maximum of 74. The null model for the model-based approach shows a significant second-level standard deviation ($\sigma_u = 0.565$; p -value < 0.0001).¹ The intraclass correlation coefficient is rather high ($\rho = 0.088$) if compared with the usual values observed

¹The reported p -value is based on the likelihood-ratio (LR) test but it should be noted that the null hypothesis for this test is on the boundary of the parameter space because it refers to a variance component. As a consequence, the LR test does not have the usual central chi-square distribution with one degree of freedom but it is better approximated as a 50:50 mixture of central chi-squares with zero and one degree of freedom [5].

for similar models, indicating the presence of a substantial association in the responses within each course of study. These preliminary results confirm that a two-level approach is suitable for analyzing our data and thus we proceed with the successive steps. After introducing the student-level covariates (Table 1), the model shows a lower but still significant second-level standard deviation ($\sigma_u = 0.332$; p -value = 0.0110) and, accordingly, a reduced value of the intraclass correlation coefficient ($\rho = 0.032$).

For this model, the value of the constant defines a reference profile, corresponding to a hypothetical student having a value of zero for each explanatory variable. This profile describes a student enrolled to the same course of registration, who has a regular or not regular registration by 1–2 years, who did not repeat years during high school, who has a school final mark less than 90, who attends lessons more or less regularly, who has relationships with other students, who does not have a job, who has no children, and who did not attend grammar school. For the reference profile the course-effect is also null. The magnitude of the between-courses unobserved heterogeneity can be investigated by computing the probabilities obtained by adding and subtracting a value of the random effect equal to twice the estimated standard deviation. These values approximately correspond to the 2.5 percentile of a normal distribution and differentiate courses with a low or a high percentage of inactive students, still referred to the reference profile. Taking into account the student-level covariates, the gap between courses with a high or a low inactivity rates is about two decimal points (0.341–0.121). For to the design-based approach (Table 2) the results show that the two approaches yield roughly similar estimates.

The Asparouhov's [14] informativeness measure (I_2) takes on a value lower than 1 for all the parameters except for the constant, whose value is 3.203, indicating a highly significant difference between the two approaches. This result implies very different values for the probability of being in the inactive status for students with a reference profile, whose probabilities are, respectively, 0.268 for the model-based approach and 0.070 for the design-based approach. In addition, the gap between courses with a high or a low inactivity rates is narrower for the design-based approach (0.107–0.028) due to a lower between-courses variance. Based on these results, we point out that since the reference profile refers to students in a relatively favorable condition, the design-based estimates seem to be closer to reality and, for this reason, preferable. Additionally, though model-based estimates have smaller standard errors, the difference with the design-based approach is not remarkable and then acceptable. In this case the sampling design is informative because a stratification variable is the response variable. The results obtained comparing model-based and design-based estimates confirm that when the sample selection probabilities depend on the values of the model response variable, even after conditioning on auxiliary variables, the sampling mechanism becomes informative and the selection effects need to be accounted for in the inference process.

Following the estimates of the design-based model and focusing on the effects of the significant student-level covariates, we can see that the probability of being in the inactive status for the reference profile (0.070) increases significantly above all for students enrolled in a not regular condition by more than 2 years (0.241)

and for students having a job (0.192). As concerns the variables characterizing university experience, the probability is higher for students not attending lessons (0.132), for students enrolled in a course different from that of registration (0.120), and for students not having contacts with classmates outside the university (0.099). Finally, about the variables characterizing the student experience before entering the university, the probability of being in the inactive status is higher for who repeated years during high school (0.109) but lower for students with a school final mark higher than 90/100 (0.047) and for students attending grammar school (0.050), respectively. Finally, motivations of registration and the status of freshmen have no significant effect. Adding course-level covariates does not produce significant effects, therefore the final model remains unchanged.

3.2 Second Cycle Courses Degree

The analysis comprises 371 students within 113 courses of study. The average number of students for each course is 3.3 with a minimum of 2 and a maximum of 26. The null model shows that the second-level standard deviation is not significant ($\sigma_u = 0.001$; p -value = 1.000). The final model (Table 3) is obtained following a design-based approach, still because of a significant difference observed in the constant. This model is much simpler from the other not only for the absence of a course-effect but also for the limited number of covariates having a significant effect on the response variable. These results confirm our hypothesis about the difference between the two types of degree programmes and support further our decision to estimate separate models. For this model, the value of the constant parameter defines a reference profile corresponding to a student who attends lessons more or less regularly, who does not have a job, who has no children, who has a regular or not regular registration by 1–2 years, and, finally, who decided to enroll mainly not for future job opportunities. The associated probability of being in the inactive status is very low (0.044). Compared to the reference profile, the probability is higher for students enrolled in a not regular condition by more than 2 years (0.235) and for students not attending lessons (0.115). As concerns the variables characterizing personal life, the probability is higher for students having a job (0.091) and having children (0.077). Finally, we introduced into the model a variable concerning the motivations of enrollment, though it is poorly significant. This variable indicates that students motivated mainly by future job opportunities have a lower probability of being in the inactive status (0.025).

4 Discussion and Final Remarks

The aim of this paper was to analyze the issue of inactive students. The literature on this topic is poor and a focus is important not only for the negative effects of inactivity on students and their families but also because an increasingly amount of the FFO rewarding share is allocated to the universities taking into account

indicators related to the students' career. A sample survey was carried out by the Statistical Observatory of the University of Pisa by selecting a stratified simple random sample of students enrolled at the University of Pisa in the academic year 2010–2011. Interviews were administered through a CATI system. Data analysis is performed using separate models for the first and the second cycle degree courses. The probability of being in the inactive status is modelled using a two-level random intercepts logistic regression model. To assess the informativeness of the sampling design, analysis was performed by comparing parameter estimates under the model-based and the design-based approaches.

The results show a considerable difference between first and second cycle degree courses on the probability of being in the inactive status. A significant course-effect is observed only for the first cycle degree even though we do not find covariates for explaining this influence. In addition, school path before entering the university (years repeated during school, final school mark, and attending grammar school) plays a crucial role for students of the first cycle but is irrelevant for students of the second cycle degree programmes. More specifically, the probability of being in the inactive status is lower for first-cycle students who have contacts with fellows outside the university and for second-cycle students whose registration is motivated by job opportunities. Finally, in common between first and second cycle are university experience (attending lessons and being in a regular or not regular registration by 1–2 years), with a positive effect and, as concern personal conditions, having a job, with a negative effect.

These results, for both first and second cycle degree, are not inconsistent with the possibility that a significant portion of inactive students are the so-called working-students. These students typically cannot attend lessons, have no relationships with other students, are enrolled for a long time, and are far from the university system. For these students, the actual university system, whose activities (lessons, consulting, and tutorship) are carried out mainly in the morning or in the early afternoon, might be hardly practicable. A simplest possible short path to control and reduce the proportion of inactive students would seem to enlarge the offer of evening and night classes. Nevertheless there is opportunity also for the improvement of counselling and services of vocational guidance, *ex ante* and *in itinere*.

References

1. Benvenuto, G.: Percorsi di studio universitari: monitoraggio delle matricole e indagine sugli studenti "inattivi". Nuova Cultura (2010)
2. Cingano, F., Cipollone, P.: University Drop-Out: The Case of Italy. Banca d'Italia, Roma (2007)
3. Comitato Nazionale di Valutazione del Sistema Universitario (CNVSU). Undicesimo Rapporto sullo Stato del Sistema Universitario. Roma (2011)
4. Raudenbush, S., Bryk, A.: Hierarchical Linear Models, 2nd edn. Sage Publications, Thousand Oaks (2002)
5. Snijders, T.A., Bosker, R.: Multilevel Analysis. An Introduction to Basic and Advanced Multilevel Modelling, 2nd edn. Sage, Thousand Oaks, CA (2011)

6. Hox, J.J.: *Multilevel Analysis: Techniques and Applications*, 2nd edn. Routledge, New York (2010)
7. Goldstein, H.: *Multilevel Statistical Models*. Wiley, New York (2011)
8. Skrondal, A., Rabe-Hesketh, S.: *Generalized Latent Variable Modelling: Multilevel, Longitudinal and Structural Equation Models*. Chapman & Hall/CRC, Boca Raton, FL (2004)
9. Särndal, C.E.: Design-based and model-based inference in survey sampling. *Scand. J. Stat.* **5**, 27–52 (1978)
10. Pfeffermann, D.: The role of sampling weights when modeling survey data. *Int. Stat. Rev.* **61**, 317–337 (1993)
11. Pfeffermann, D., Sverchkov, M.: Fitting generalized linear models under informative probability sampling. In: Skinner, C., Chambers, R. (eds.) *Analysis of Survey Data*, Wiley, New York, pp. 175–195 (2003)
12. Pfeffermann, D., Skinner, C.I., Holmes, D.I., Goldstein, H., Rasbash, I.: Weighting for unequal selection probabilities in multi-level models. *J. R. Stat. Soc.* **60**, 23–56 (1998)
13. Grilli, L., Pratesi, M.: Weighted estimation in multi-level ordinal and binary models in the presence of informative sampling designs. *Surv. Methodol.* **30**, 93–103 (2004)
14. Asparouhov, T.: General multilevel modeling with sampling weights. *Commun. Stat. Theory Methods* **35**, 439–460 (2006)
15. Rabe-Hesketh, S., Skrondal, A.: Multilevel modelling of complex survey data. *J. R. Stat. Soc.* **169**, 805–827 (2006)

Part IV

Social Statistics, Demography and Health Data

The Material Deprivation of Foreigners: Measurement and Determinants

Annalisa Busetta, Anna Maria Milito, and Antonino Mario Oliveri

Abstract

We examine the material deprivation of foreigners on a sub-sample of the 2009 Italian Survey on Income and Living Conditions carried out by Istat. We employ an index of material deprivation that takes into account the regional level of analysis, and relies on the assignment of weights to deprivation items. The effects produced on material deprivation by several variables, interpreted as determinants, are investigated through a zero-inflated beta regression model.

1 Introduction

Although differing in methods and range of controls used, statistics on living conditions (in terms of poverty, social exclusion, deprivation) point foreigners as one of the more vulnerable groups [1–3]. Despite these evidences, up to the present there are only few studies in Italy on living conditions of foreigners, mainly based on ad hoc sample surveys.¹

¹The first official sample surveys on foreigners are very recent. We refer to the Istat IT-Silc 2009 survey on households with foreigner members (recently released) and to the 2008 Labour force survey on the integration of immigrants.

A. Busetta

Department of Economics, Business and Statistics, University of Palermo, Palermo, Italy
e-mail: annalisa.busetta@unipa.it

A.M. Milito (✉) • A.M. Oliveri

Department of Cultures and Society, University of Palermo, Palermo, Italy
e-mail: annamaria.milito@unipa.it; antoninomario.oliveri@unipa.it

Material deprivation has often been considered as an indirect measure of permanently low income [4] and a valid indicator of multidimensional poverty [5, 6]. It is defined as "... inability for individuals or households to afford those consumption goods and activities that are typical in a society at a given point in time, irrespective of people's preferences with respect to these items" [7]. Moreover most authors define material deprivation as "exclusion from the minimum acceptable way of life in one's own society because of inadequate resources" [8–13]. Another common definition refers to "the lack of socially perceived necessities" [11, 14]. Therefore studying the material deprivation of foreigners provides us also with information on social inclusion and living standards [15].

In this paper we investigate the condition of material deprivation experienced by foreigners in Italy through a weighted version of the official Eurostat index of material deprivation, in which the weights take into account the relevance of each deprivation item at the regional level. This version of the index allows to compare individuals in terms of diffusion and intensity of material deprivation. This index is used as a response variable within a regression model fitted to estimate the impact of some key explanatory variables on material deprivation. Data are a sub-sample of foreigners drawn from the 2009 Survey on Income and Living Conditions for Italy (henceforth IT-Silc). The sample consists of 1633 foreigners.

2 Literature Background

Much of the interest in measuring material deprivation stems from the work of Townsend [16], who correlated the concept of deprivation to the broader notion of ability to enjoy an acceptable standard of living. In his seminal paper Townsend measured 11 forms of deprivation through a set of 60 indicators with binary deprivation scores (i.e. having or not having a specific good), which were finally synthesised in a composite indicator. Subsequent contributions have both criticised and extended Townsend's measurement approach. In particular Piachaud [17] questioned the failure to distinguish between the lack of a good (or an activity) due to a voluntary choice of individuals and the lack resulting from financial constraints. Ringen [18] criticised Townsend's approach complaining that it assesses material deprivation (intended as a direct measure of poverty) through an income threshold (an indirect measure of poverty). Other authors have raised questions on the arbitrary list of items used and on the failure to take into account the seriousness of different forms of deprivation [19].

By this time the literature on material deprivation has overcome most of these criticisms and converges on some widely accepted characteristics (for a review see [5]):

- the household is the fundamental unit within which resources are shared and needs are satisfied;

- individuals who cannot afford a certain good or service must be distinguished from those who do not have this good or service for different reasons (including free personal choice);
- the items selected should measure differences in deprivation rather than differences in tastes and preferences;
- material deprivation measures are expected to be consistent with both absolute and relative interpretations of poverty.

Nowadays the overall measures of material deprivation are based on multiple binary indicators referring to whether households lack various items and activities that are perceived as necessities and their lack is because they cannot afford them rather than because they do not like them, i.e. an “enforced lack” ([16] and later [19–21]).

Also in the studies on separate dimensions of deprivation, multiple binary indicators are combined into a single numerical scale [10, 11, 22–25].

Eurostat released an “official” operational definition of material deprivation, from which a composite indicator was proposed, made up of nine binary elementary indicators [26]. Nowadays the Eurostat’s solution is the most frequently cited and we ourselves recurred to it in this study.

The debate on item selection still draws the attention of scholars. In particular, alternative operational definitions have been proposed in order to measure material deprivation, which may determine the selection of different sets of items ([5, 15] and references therein; [27]). None of them proved to be free from drawbacks yet, and this is why we decided to adopt in this study the official Eurostat definition, which features two relevant advantages: (1) it allows easy cross-country comparisons; (2) it follows most of the desirable features of measurement scales [28]. Eurostat is currently revising its measurement system in order to propose a new set of more reliable items [29, 30].

3 Descriptives

It is a matter of fact that in Italy material hardship is higher among foreigners both in terms of diffusion and intensity [3]: in 2009 74.4 % of foreigners experienced an enforced lack of at least one item (less than 50 % among natives); around 34 % were deprived of at least 3 out of 9 items (14 % among natives) and around 20 % of at least 4 items (6.1 % among natives) (see Table 1).

Moreover Southern regions appear in general to be poorer and more deprived than the Northern ones, and this is also true for deprivation of foreigners within regions: the percentage of foreigners experiencing deprivation in respect to at least one item varies from 0 % of Molise and Sardinia to around 60 % for the autonomous province of Trento (see Table 2). The incidence of material deprivation (i.e. having at least 3 out of the 9 deprivation items, according to the Eurostat official definition)

Table 1 Incidence and intensity of material deprivation in 2009

	Not deprived	At least 1 deprivation	Incidence of		Intensity of deprivation ^a	
			Deprivation ^b	Hard deprivation ^c	Mean	Standard deviation
Foreigners	25.6	74.4	33.8	18.6	2.6	1.4
Natives	51.8	48.2	14.0	6.1	2.1	1.2
All	50.3	49.7	15.2	6.8	2.1	1.2

Source: IT-Sile survey 2009

^aOnly among people who have at least 1 deprivation

^bWith at least 3 out of 9 items

^cWith at least 4 out of 9 items

Table 2 Incidence and intensity of material deprivation among foreigners in 2009

	Not deprived	At least 1 deprivation	Incidence of		Intensity of deprivation ^a		
	% (No. of households)		Deprivation ^b	Hard deprivation ^c	Mean	Standard deviation	Obs
Piedmont	21.3(25)	78.7	28.1	20.8	2.4	1.4	81
Aosta Valley	46.4(9)	53.6	5.5	2.6	1.5	0.8	17
Lombardy	21.6(64)	78.4	35.3	19.9	2.6	1.3	226
Bozen-Bolzano	27.4(18)	72.6	30.3	14.8	2.1	1.2	31
Trento	59.1(11)	40.9	17.2	4.6	2.3	1.6	10
Veneto	28.6(42)	71.5	30.7	17.9	2.7	1.7	122
Friuli-Ven. Giulia	15.9(15)	84.1	24.3	14.6	2.3	1.0	63
Liguria	27.7(26)	72.3	25.2	2.4	2.0	0.9	39
Emilia-Romagna	30.9(44)	69.1	27.0	9.8	2.2	1.2	104
Tuscany	35.7(47)	64.3	18.2	8.2	2.3	0.9	90
Umbria	26.5(38)	73.5	46.1	22.9	2.9	1.0	77
Marche	30.4(17)	69.7	36.4	26.8	2.8	1.2	58
Lazio	36.5(57)	63.5	28.9	11.7	2.7	1.4	91
Abruzzo	6.7(2)	93.3	38.7	10.6	2.1	1.1	40
Molise	0.0(0)	100.0	48.6	0.0	2.4	0.7	8
Campania	12.9(3)	87.1	49.7	34.6	3.1	1.7	29
Apulia	14.0(7)	86.0	40.3	31.0	2.7	1.5	31
Basilicata	8.2(2)	91.8	81.8	71.9	3.6	1.0	16
Calabria	13.5(4)	86.5	68.1	25.7	3.2	1.5	25
Sicily	9.9(7)	90.1	75.2	63.6	4.0	1.5	33
Sardinia	0.0(0)	100.0	57.6	3.3	2.2	1.3	4

Source: Sub-sample of IT-Sile survey 2009

^aOnly among people who have at least 1 deprivation

^bWith at least 3 out of 9 items

^cWith at least 4 out of 9 items

varies from 5.5 % of the Aosta Valley to 81.8 % of Basilicata and the intensity of deprivation (mean number of deprivation items) varies from 1.5 of the Aosta Valley to four of Sicily.

According to the last fiscal federal reform, Italian municipalities and regions are charged with the implementation of the most part of social policies. These are the reasons why scholars have recently adopted the regional level to analyse material deprivation in Italy [31].

4 The Measurement of Material Deprivation

As reported above, we measured material deprivation as “the number of individual’s (enforced) lack of access to essentials of life” [5]. Relying on Eurostat’s definition [26] and to data available from the IT-Silc survey, nine “essential” items were considered: (1) to face unexpected expenses; (2) 1 week annual holiday away from home; (3) to pay for arrears; (4) a meal with meat, chicken or fish every second day; (5) to keep home adequately warm; (6) a washing machine; (7) a colour TV; (8) a telephone, (9) a personal car.

The most used deprivation scales assume that each deprivation item has the same importance. Nevertheless, the issue of item weighting has been broadly considered in the literature² and many solutions have been suggested [5, 31].

We calculate the Material Deprivation (MD) index of foreigners in Italy as:

$$MD_i = \sum_{j=1}^9 w_{jk} X_{ij} \tag{1}$$

where $i = 1, 2, \dots, N$ indicate the sample units, $j = 1, 2, \dots, 9$ the deprivation items, $k = 1, 2, \dots, 21$ are the 19 Italian regions and 2 autonomous provinces (henceforth all simply named “regions”), X_j represent indicator deprivation variables (1 = not owning, 0 = owning the item) and the w_{jk} terms represent normalised weights³ calculated as follows:

$$w_{jk} = \frac{h_{jk}}{\sum_{j=1}^9 h_{jk}} \text{ and } \forall js \sum_{j=1}^9 w_{jk} = 1 \tag{2}$$

Due to the regional nature of deprivation in Italy, it seems reasonable that such weights may vary across regions. In Eq. (2), h_{jk} is the weight calculated for the j -th item and the k -th region on the whole IT-Silc sample. It represents the percentage of sample units *having* the item so that higher weights are given to more possessed items in the region; that means, individuals result more deprived if they *do not have* an item possessed by most of the people living in the same region. These

²For a detailed review on the issue see the recent contribution by Decancq et al. [35].

³For a review of alternative weighting techniques see Nardo et al. [36].

weights, usually called “prevalence weights”, are consistent with the concept of material deprivation as a relative phenomenon, which we stressed in the above. The approach is opposed to the “consensus weights” one, which is based on the proportion of people considering the item i as “absolutely necessary or necessary” in the Eurobarometer survey.

The MD_i index measures the intensity of deprivation through the weighted mean of items not owned by the individuals in our sample. It is expressed as a continuous variable in the $(0, 1)$ interval; the lower bound (zero) is achieved by the individual when he/she does not suffer from any deprivation, whereas the upper bound (one) corresponds to the lack of all considered items.

Finally it is not redundant to say that the choice of sum as the aggregation function is chosen as material deprivation items are assumed to be mutually independent, even if they measure the same latent trait. The validity of the sum-score approach is supported by the recent research output by Jenkins and Cappellari [32]. They maintain the weakness of the theoretical foundations of the ‘sum-score’ scales, showing that results obtained through item response modelling approach provided “similar pictures of deprivation patterns and their determinants, and so our results might be construed as providing an empirical rationale for the sum-score approach”.

5 Methodology

Under strong regional inequalities in material deprivation (see Sect. 3), a major aim of this paper is to investigate the determinants of deprivation among foreigners in Italy. This was done by interpreting MD as a response variable and some socio-demographic characteristics of foreigners and of their households as explanatory variables.⁴

The histogram and the box-plot in Fig. 1 present the distribution of foreigners by their scores on the individual deprivation index. The clump-at-zero (the bar with the dot above) in the histogram represents 26.8 % of foreigners who are not deprived in any of the nine dimensions. We observe that the distribution of the data is asymmetric, right-skewed with an inverted “J” shape (the mean equals to 0.18, median 0.17 and st. dev. 0.16; trimmed distribution: $n = 1195$, mean 0.24, median 0.20, st. dv. 0.14). Among foreigners no one lives in households with the maximum intensity of deprivation.

Respect to the choice of a proper statistical model, as the MD index assumes a non-negligible number of zeros, we referred to mixed continuous–discrete distributions. This family of distributions, introduced by Ospina and Ferrari [33] and usually defined zero-or/and-one inflated beta distributions, allows to model data that assume values in $[0, 1)$, $(0, 1]$ or $[0, 1]$. In order to model the MD distribution via this model, we had to evaluate if the observed inflation of zeros was generated from

⁴Explanatory variables were chosen recurring to the stepwise procedure and suggestions from the literature on living conditions and on material deprivation of foreigners in Italy.

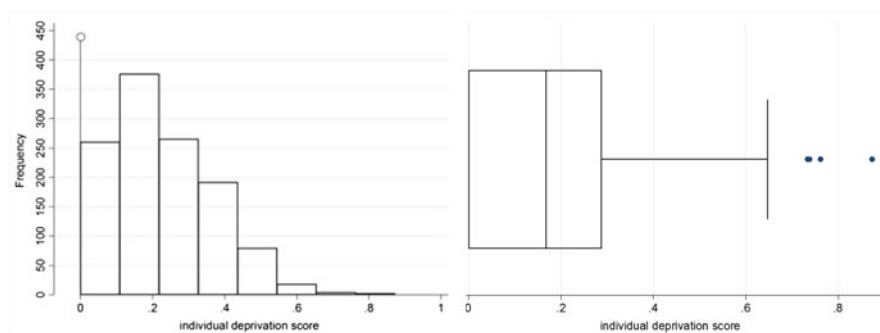


Fig. 1 Frequency histogram and box-plot of the individual deprivation (MD) index

a distinct process or occurred through the same process as the other values (if this last was the case, a fractional logit model would be more appropriate). For material deprivation, the processes involved were (1) that an individual was deprived vs. not deprived. If not deprived, the only possible outcome was zero; (2) if deprived, it was then a process with varying intensity. The expected value was expressed as a combination of the two processes. Therefore, it seemed adequate to use a zero-inflated beta (ZOIB) model with two components: (1) a logistic regression model for whether or not individuals have no deprivation ($MD = 0$) and (2) a beta model for a degree of deprivation between 0 and 1, which are simultaneously estimated.

6 Model and Discussion

The zero-inflated beta model estimates the impact of socio-economic and demographic characteristics⁵ of foreigners and their households, controlling for the regions where they lived. The model fits data well. Estimates⁶ show that the probability of experiencing material deprivation depends on a range of characteristics of foreigners and of the households they live in (see Table 3). The relation between

⁵The explanatory variables included in the zero-inflated beta model are: gender, age, education (up to upper secondary school, higher secondary school, tertiary school), labour market position (employee, self-employed, unemployed, inactive), household composition (one person household, two adults without children, single parent household with children, two adults with children, other household with dependent children), working intensity (WI) status in four classes as defined by Eurostat, having an Italian partner (respect to a foreign partner), tenants (respect to home owners), EU citizen (respect to extra-EU), self-assessed health (bad or very bad).

⁶In the beta model positive estimates indicate the amount of increase in the MD index that would be due to an increase (or to a change in state) in the explanatory variables, whereas in the logit model positive estimates indicate the amount of increase in the predicted probability of being not deprived. This is the reason why in most cases the same covariates show opposite signs in the two parts of the zero-inflated beta model. All the estimates are reported on the logit scale.

Table 3 Zero-inflated beta model^a

	Beta component		Logit component	
	Coeff.	Std. err.	Coeff.	Std. err.
<i>Personal characteristics</i>				
Woman (ref. man)	-0.041	0.043	0.057	0.138
Age	0.027	0.009	-0.061	0.026
Squared age	0.000	0.000	0.001	0.000
<i>Education (ref. up to upper secondary school)</i>				
Higher Secondary School	-0.163	0.043	0.226	0.141
Tertiary School	-0.171	0.082	0.917	0.206
<i>Labour market position (ref. employee)</i>				
Self-employed	-0.082	0.086	0.591	0.226
Unemployed	0.274	0.063	-0.928	0.262
Inactive	0.021	0.057	-0.078	0.182
<i>Household composition (ref. one person household)</i>				
Two adults without children	0.009	0.076	0.657	0.231
Single parent household with children	-0.057	0.076	0.123	0.239
Two adults with children	0.105	0.124	-0.056	0.422
Other household with dependent children	-0.163	0.086	0.532	0.265
<i>Working intensity status (WI) (ref. WI = 0)</i>				
0 < WI < 0.5	-0.221	0.109	-0.591	0.354
0.5 ≤ WI < 1	-0.319	0.106	-0.306	0.343
WI = 1	-0.421	0.110	-0.051	0.357
<i>Other characteristics</i>				
Having an Italian partner (ref. a foreign partner)	-0.152	0.067	0.728	0.173
Tenants (ref. owners)	0.203	0.046	-0.553	0.132
EU citizen (ref. extra-EU citizen)	-0.087	0.048	0.386	0.138
Bad or very bad self-assessed health (ref. not bad/very bad)	0.070	0.094	-1.157	0.403

Parameter estimates (coeff.) and standard errors (std. err.)

^aControlled for regions

age and deprivation shows that young foreigners are more exposed to the risk of deprivation and this risk declines with age.

Both education and labour market position have an important role in reducing the intensity of deprivation, and increase the probability of being no deprived. In particular, having a tertiary educational level (compared to low educational level) is associated with a lower probability of being no deprived and with lower intensity of deprivation. Also being unemployed (vs. being employed) acts in the expected direction since it is associated with low chance of being no deprived and high intensity of deprivation.

As expected, the work intensity status of the household⁷ has an important role [34]: the higher the share of workers in the household, the higher the probability of experiencing more severe forms of deprivation. The work intensity status seems to have no effect on the odds ratio of being deprived.

In respect to the household structure, our estimates reveal that having a native partner is important both for protecting foreigners from falling into deprivation and for reducing the intensity of deprivation. Also the household composition has a role in protecting from deprivation: foreign couples without children and large households with dependent children are on average less deprived than other households (ref. single households). The household composition seems to have no effect on the intensity of deprivation yet.

Studies on deprivation agree that material deprivation affects far more sick and disabled people than the rest of the population [5]. Our analysis confirms that individuals who self-assess bad or very bad health are more likely to experience deprivation.

Home tenure also helps to understand material deprivation. Home owners are less likely to fall into material deprivation than renters and even if deprived, they have a less intense experience.

7 Conclusions and Limitations

This paper sheds light on the new field of material deprivation of foreigners in Italy, taking into account the regional nature of deprivation and explicitly including item weighting. Nevertheless, the results of the model suffer from weakness of external validity because data on foreigners are drawn from the standard IT-Silc survey whose frame population is not that of foreigners but that of all Italian households.

We are conscious that the years since migration and the age at arrival have a direct effect on the ability to enjoy a standard of living that is generally considered acceptable, and an indirect effect on all the explanatory variables included in our model. Unfortunately, these variables are not included in the survey and they were consequently omitted from our model. All these limitations will be overcome working on the sample of 6000 households included in the ad hoc survey on foreigners carried out by Istat in 2009. A further development of this work will consist in the evaluation of the differences in living conditions among different nationalities in Italy.

Acknowledgements Although the paper is common responsibility of all authors, Sects. 1 and 7 can be attributed to A.M. Milito, Sects. 2 and 4 to A.M. Oliveri and Sects. 3, 5 and 6 to A. Busetta.

⁷The work intensity status, computed at the household level, summarises the work status over the past year for all work age household members (from 18 to 64 years old). This measure is calculated by Eurostat taking into account the ratio of worked months over workable months, averaged over all work age household members and categorised as follows: $WI = 0$, $0 < WI < 0.5$, $0.5 \leq WI < 1$, $WI = 1$.

The paper has received financial support from the University of Palermo [grant no. ORPA07ZMAE under the responsibility of Anna Maria Milito].

References

1. Caritas Italiana - Fondazione E. Zancan: Rapporto 2011 su povertà ed esclusione sociale in Italia. Il Mulino, Bologna (2011)
2. Commissione di Indagine sull'Esclusione Sociale (CIES): Rapporto sulle politiche contro la povertà e l'esclusione sociale. CIES, Roma (2011)
3. Istat: Households with foreigners: indicators of economic distress. Notes for the press - 28 February 2011. Istat, Rome (2011)
4. Willits, M.: Measuring child poverty using material deprivation: possible approaches. Working paper of Department of work and pensions, n. 28 (2006). <http://research.dwp.gov.uk/asd/asd5/WP28.pdf>. Cited 11 February 2012
5. Guio, A.C.: What can be learned from deprivation indicators in Europe?. Eurostat Meeting Working Paper, pp. 1–33 (2009)
6. Nolan, B., Whelan, C.T.: Multidimensionality of poverty and social exclusion. In: Jenkins, S.P., Micklewright, J. (eds.) Resources, Deprivation and Poverty. Clarendon Press, Oxford (2007)
7. OECD: Material deprivation. In: Glossary of statistical terms. OECD (2007). <http://stats.oecd.org/glossary/detail.asp?ID=7326>. Cited 11 February 2012
8. Callan, T., Nolan, B., Whelan, C.T.: Resources, deprivation and the measurement of poverty. *J. Soc. Policy* **22**(2), 141–172 (1993)
9. Kangas, O., Ritakallio, V.: Different methods – different results? Approaches to multidimensional poverty. In: Andress, H.-J. (ed.) Empirical Poverty Research in a Comparative Perspective. Ashgate Publishing Limited, Aldershot (1998)
10. Layte, R., Nolan, B., Whelan, C.T.: Reassessing income and deprivation approaches to the measurement of poverty in the Republic of Ireland. *Econ. Soc. Rev.* **32**(3), 239–261 (2001)
11. Nolan, B., Whelan, C.T.: Measuring poverty using income and deprivation indicators: alternative approaches. *J. Eur. Soc. Policy* **6**(3), 225–240 (1996)
12. Perry, B.: The mismatch between income measures and direct outcome measures of poverty. *Soc. Policy J. New Zealand* **19**, 101–127 (2002)
13. Whelan, C.T., Layte, R., Maître, B.: Multiple deprivation and persistent poverty in the European Union. *J. Eur. Soc. Policy* **12**(2), 91–105 (2002)
14. Bradshaw, J., Finch, N.: Overlaps in dimensions of poverty. *J. Soc. Policy* **32**(4), 513–525 (2003)
15. Whelan, C.T., Nolan B., Maître, B.: Measuring material deprivation in the enlarged EU. ESRI working paper, n. 249 (2008). http://irserver.ucd.ie/dspace/bitstream/10197/1014/1/nolanb_workpap_028.pdf. Cited 11 February 2012
16. Townsend, P.: Poverty in the United Kingdom. Penguin Books, Harmondsworth (1979)
17. Piachaud, D.: Peter Townsend and the Holy Grail, New Society, (1981)
18. Ringen, S.: Direct and indirect measures of poverty. *J. Soc. Policy* **17**, 351–65 (1988)
19. Gordon, D., Adelman, L., Ashworth, K., Bradshaw, J., Levitas, R., Middleton, S., Pantazis, C., Patsios, D., Payne, S., Townsend, P., Williams, J.: Poverty and Social Exclusion in Britain. Joseph Rowntree Foundation, York, pp. 101 (2000)
20. Gordon, D., Pantazis, C. (eds.): Breadline Britain in the 1990s. Ashgate Publishing Limited, Aldershot (1997)
21. Mack, J., Lansley, S.: Poor Britain. George Allen and Unwin, London (1985)
22. Atkinson, A.B.: Multidimensional deprivation: contrasting social welfare and counting approaches. *J. Econ. Inequal* **1**, 51–65 (2003)
23. Layte, R., Maître, B., Nolan, B., Whelan, C.T.: Persistent and consistent poverty: an analysis of the first two waves of the European Community Household Panel. *Rev. Income Wealth* **47**, 427–440 (2001)

24. Nolan, B., Whelan, C.T.: Resources, Deprivation and Poverty. Oxford University Press, Oxford (1996)
25. Whelan, C.T., Layte, R., Maitre, B., Nolan, B.: Income, deprivation and economic strain: an analysis of the European Community Household Panel. *Eur. Sociol. Rev.* **17**(4), 357–372 (2001)
26. Eurostat: Material deprivation. In: Eurostat “Living conditions and social protection glossary”. http://epp.eurostat.ec.europa.eu/statistics_explained/index.php/Glossary:Material_deprivation. Cited 22 November 2013
27. Fusco A., Guio A. C., Marlier E.: Building a material deprivation index in a multinational context: lessons from the EU experience. In: Berenger, V., Bresson, F. (eds.) *Poverty and Social Exclusion around the Mediterranean Sea Economic Studies in Inequality. Social Exclusion and Well-Being*, pp. 43–71. Springer US, New York (2013)
28. Guio, A.C.: Material Deprivation in the EU, *Statistics in Focus, Population and Social Conditions, Living Conditions and Welfare*. Eurostat, Luxembourg (2005)
29. Eurostat: 2009 EU-SILC Module on Material Deprivation. Assessment of the Implementation. Eurostat, Luxembourg (2011). http://epp.eurostat.ec.europa.eu/portal/page/portal/income_social_inclusion_living_conditions/documents/tab2/Analysis%202009%20Module.pdf. Cited 22 November 2013
30. Eurostat: Measuring Material Deprivation in the EU. Indicators for the Whole Population and Child-Specific Indicators. Eurostat, Luxembourg (2011)
31. D’Ambrosio, C., Giuliano, G., Tenaglia, S.: Material deprivation: an application to Italian regions. *Polit. Econ.* **3**, 349–368 (2009)
32. Jenkins, S.P., Cappellari, L.: Summarizing multiple deprivation indicators. In: Micklewright, J. (ed.) *Inequality and Poverty Re-Examined*, pp. 166–184. Oxford University Press, Oxford (2007)
33. Ospina, R., Ferrari, S.L.P.: Inflated beta distributions. *Stat. Pap.* **51**, 111–126 (2010)
34. Eurostat: Income, Poverty and Social Exclusion: Second Report. Eurostat, Luxembourg (2002)
35. Decancq, K., Lugo, M.A.: Weights in multidimensional indices of well-being: an overview. *Discussions Paper Series (DPS) of Center for Economic Studies - Katholieke Universiteit Leuven 10.06* (2010) Available via: <http://www.econ.kuleuven.be/research/CES/>. Cited 10 October 2012
36. Nardo, M., Saisana, M., Saltelli, A., Tarantola, S.: *Tools for Composite Indicators Building*. European Commission, Ispra (2011)

How Do Life Course Events Affect Paid and Unpaid Work of Italian Couples?

Maria Gabriella Campolo, Antonino Di Pino, and Ester Lucia Rizzi

Abstract

The paper analyzes the impact of life course events, and in particular of parenthood, on the paid and unpaid working activity of dual-earner couples in Italy. To this purpose, we use the panel dataset provided by the 2003–2007 Istat Multipurpose Survey. To correct misspecification due to unobserved variables, we adopt a difference-in-differences specification of simultaneous equations of market and domestic work supply. Our results show that the negative effect of transition to parenthood on female paid work supply is stronger than the positive effect of wages.

1 Introduction

Several studies analyse empirically the paid and unpaid work supply of women and their partners adopting a life course perspective [1, 13]. However, the estimated effect of life course events (such as parenthood) may be inconsistent as a consequence of the misspecification of latent factors involving unobserved effects over time and the unobserved heterogeneity “between” individuals [2, 3]. In this study, we aim to solve these problems by adopting a partially original methodological approach that allows us to correct for the effects of omitted variables. To this end, a correction method is here applied to a longitudinal simultaneous equations model

M.G. Campolo • A. Di Pino (✉)

Dipartimento di Economia, University of Messina, Via T. Cannizzaro, 278, 98122 Messina, Italy
e-mail: mgcampolo@unime.it; dipino@unime.it

E.L. Rizzi

Université Catholique de Louvain, Louvain-la-Neuve, Belgium
e-mail: ester.rizzi@uclouvain.be

of the paid and unpaid work of both partners. Misspecification effects are corrected by imposing non-null covariances in the error terms between repeated observations over time and across the equations. In this way, the parenthood transition effect on partners' paid and unpaid work can be evaluated simultaneously, taking into account that partners' decisions on the allocation of working time are mutually correlated. In addition, the life course events effects are specified using a two-wave Difference-in-Differences (*DID*) approach. This approach allows us to disentangle the actual life event effect from the selective nature of those who experience the life course transition during the period of survey.

For the empirical analysis, we use a sample of 562 two-earner couples, interviewed in 2003 and 2007. The source is the Istat (Italian National Institute of Statistics) Multipurpose Panel Survey 2003–2007.

Our estimation results show that the negative impact of motherhood on women's labour supply is stronger than the positive impact of wages. Moreover, hourly wage negatively affects female domestic work. This result may be interpreted as an empirical confirmation that potential hourly wage represents for women a measure of the opportunity cost on unpaid work and informal care activity [4, 6]. Evaluating the impact of gender attitudes on the allocation of time, we found that the more traditional couples adopt a more gendered division of both paid and unpaid work after childbearing, independently of the specific earnings ability of each partner.¹ The paper is organized as follows: in the next section we explain the rationale of our methodology in relation to recent developments in life course analysis. In Sect. 3, we report the characteristics of the dataset and the model specification. In Sect. 4 we discuss our results.

2 Methodology

In this paper we consider two misspecification problems that may affect the estimation results of a life course transition model: (1) the influence of latent factors determining change over time, and (2) the effect of “between individuals” omitted factors correlated with covariates (endogeneity).

Regarding the first problem, as the data include repeated measurements of the same subjects, observations over time may be correlated because latent factors that predispose the subjects to self-report their paid and unpaid work hours in a particular way initially are likely to encourage similar responses over time. The second misspecification problem in the intra-household allocation of working time between partners may arise from difficulties in specifying the bargaining process, determined by latent psychological and cultural factors, such as gender attitudes.

¹Several studies referring to different countries found that the birth of a child increases the propensity to adopt a gendered division of labour in the family [3, 12, i.a.]. However, longitudinal studies for Italy on intra-household time allocation and life course events influence are still rare.

Researchers try to identify these latent effects by introducing random effects on intercept and slopes in the model, that is, by adopting *Fixed Effect (FE)* or *Random Effect (RE)* estimators [15]. In particular, to correct for the endogeneity effects given by unobserved heterogeneity “between” individuals, *cluster-specific fixed effect* or *cluster-means model* corrected for measurement errors [7] are properly used. With regard to the latent influence of gender attitudes, several authors use as a proxy of individual gender attitudes, or as a clustering classification criterion for couples, the extent to which the subject agrees with specific statements regarding the role of women and men in the family [3, i.a.].

However, the use of proxy variables is often not sufficient to correct for endogeneity due to unobserved gender attitudes [3]. In addition, *FE* models do not provide estimation of coefficients for time-invariant variables, but only for variables changing in time [9, 10]. Finally, for the phenomenon we are focusing on, we also need a method that allows us to estimate paid and unpaid work hours of both partners simultaneously, especially taking into account the latent bargaining process between partners in working-time allocation. For this reason, we suggest, as a possible remedy to correct misspecification effects on estimates, to implement a dynamic specification of paid and unpaid work equations of both partners in a Seemingly Unrelated Regression Equations (*SURE*) model (four simultaneous equations in total), where misspecification problems given by unobserved heterogeneity between partners can be managed by imposing specific constraints on the covariance matrix [14].

More in detail, to correct the latent influence of unobserved variables, we assume that common latent factors influencing the paid and unpaid work of both partners are included in the error terms of each equation. In other terms, in the *SURE* specification the dependent variable of each equation is correlated with the dependent variables of the other equations through the unobserved components only.² Thus, non-null covariances in the repeated observation over time and non-null covariances between the error terms of each equation can be identified and used to correct estimates, respectively, for the over-time bias effect, and for the cross-sectional unobserved heterogeneity. When adopting a *SURE* specification, a Generalized Least Squares (*GLS*) procedure to obtain efficient estimates can be applied.

In addition, a *DID* specification of each equation allows us to identify the effect of an individual’s life course events. Namely, we adopt a *DID* parametrization combining a time-invariant dummy variable measuring the change of status of the subject (“status” dummy) with a dummy indicating the wave. In this way, we obtain a time-variant dummy variable indicating if the life course event has been experienced by the subject before the second wave. The estimated coefficient of this dummy variable can be considered as an unbiased measure of the impact of the life course event experienced by the subject on his/her working activity.

²Consequently, no instrumental variables are necessary to estimate the model.

3 Data and Model

We use data from the Istat Multipurpose Panel Survey, which refers to a sample of 9997 individuals. Interviews were conducted at two waves, in 2003 and 2007, and information on both paid and domestic work, life course events, and fertility was collected. The survey produced a two-wave balanced panel sample, as all individuals are surveyed over time at the time $t = 0$ (first wave), and at the time $t = 1$ (second wave). For our study we select a sub-sample of 562 Italian “two-earner” couples (1124 subjects).³ In our sub-sample women are aged 18–45 and their partners are aged 18–60 in 2003. A descriptive analysis (not presented here) shows, as expected, that men work more in the labour market than women and, conversely, women spend more time in domestic work than men, but we observe a reduction of this “gap” in the year 2007. The more children there are, the fewer hours women spend in the labour market and the more hours they spend in unpaid domestic work.

The Istat Multipurpose dataset contains no information on income. For this reason, we had to apply a matching procedure in order to import information on hourly labour income data (in Euros) from the Bank of Italy Surveys on Household Income and Wealth conducted in the years 2002 and 2006. For both years income is converted according to the 2007 price-index. Our matching procedure is based on the estimated conditional probability of assignment (propensity score) used to import information on individual wages from both the Bank of Italy surveys in the year 2002 (3640 individuals) and in the year 2006 (3242 individuals). Only individuals who live in two-earner couples are selected from the Bank of Italy dataset. To estimate the probability of assignment we use a *probit* regression, where the dependent variable is a binary dummy equal to one if the observation belongs to the Istat Multipurpose dataset, and equal to zero if the observation belongs to the Bank of Italy survey. The variables conditioning the assignment (*probit* regressors) are: gender (dummy), a dummy indicating if the subject is head of the household, the age, the geographical area of residence, the parents’ education, the working experience in years, a dummy to indicate if the subject is a house owner, the number of family members. Statistics on covariate balancing demonstrate that the matching procedure here adopted (applying a kernel-type algorithm to the estimated propensity score intervals) increases the similarity between the Istat sample and Bank of Italy sample. In order to avoid the risk of importing data characterized by self-selection, a Mantel–Haenszel test-statistics has been computed to verify if latent factors affect the assignment to treatment.⁴

The potential effects of life course events on partners’ working activity are here modelled by the specification of a simultaneous labour supply and domestic work

³We take into account only partners living together in both periods (2003 and 2007).

⁴Matching algorithm is implemented in STATA 11. More details regarding matching procedure, not presented here for the sake of brevity, are available on request.

equations model for both partners. Therefore, our model is specified as follows:

$$\ln L_{wti} = \mathbf{s}'_i \boldsymbol{\alpha}_{Lw} + \lambda_{Lw} t + t \cdot \mathbf{s}'_i \boldsymbol{\delta}_{Lw} + \mathbf{x}'_i \boldsymbol{\beta}_{Lw} + \mathbf{z}'_{it} \boldsymbol{\gamma}_{Lw} + u_{Lwti} \tag{1}$$

$$\ln D_{wti} = \mathbf{s}'_i \boldsymbol{\alpha}_{Dw} + \lambda_{Dw} t + t \cdot \mathbf{s}'_i \boldsymbol{\delta}_{Dw} + \mathbf{x}'_i \boldsymbol{\beta}_{Dw} + \mathbf{z}'_{it} \boldsymbol{\gamma}_{Dw} + u_{Dwti} \tag{2}$$

$$\ln L_{mti} = \mathbf{s}'_i \boldsymbol{\alpha}_{Lm} + \lambda_{Lm} t + t \cdot \mathbf{s}'_i \boldsymbol{\delta}_{Lm} + \mathbf{x}'_i \boldsymbol{\beta}_{Lm} + \mathbf{z}'_{it} \boldsymbol{\gamma}_{Lm} + u_{Lmti} \tag{3}$$

$$\ln D_{mti} = \mathbf{s}'_i \boldsymbol{\alpha}_{Dm} + \lambda_{Dm} t + t \cdot \mathbf{s}'_i \boldsymbol{\delta}_{Dm} + \mathbf{x}'_i \boldsymbol{\beta}_{Dm} + \mathbf{z}'_{it} \boldsymbol{\gamma}_{Dm} + u_{Dmti} \tag{4}$$

Dependent variables are given by the logarithm of weekly working hours spent in paid and domestic work by women (Eqs. 1 and 2) and men (Eqs. 3 and 4). The indexes *w* and *m* refer, respectively, to women and men. The indexes *i* and *t* refer, respectively, to the *i*-th individual and to time *t*. Dependent variables are observed on the same subject at the first time period, *t* = 0, and at the second time period, *t* = 1.

Regarding the specification of each equation, $\boldsymbol{\alpha}$ is a vector of coefficients measuring the impact of the life course events dummies included in the row vector $\mathbf{s}' = [s_1 \ s_2 \ \dots]$. Each (time-invariant) dummy signals the status of the subject as a consequence of a specific life course event. The dummy is time invariant, meaning that if the subject experiences the event in the time interval, its value is indicated as 1 at the two waves. Life course events considered in our study are the following: (1) transition to parenthood, (2) change of union status from cohabitation to marriage, (3) modification of working condition from full time to part time (and vice versa), (4) help received (purchased or informal) in domestic activity and/or in childcare, (5) dummies regarding the individual attitudes on woman’s role in the family (gender role attitudes). Consequently, $t \cdot \mathbf{s}' = [t \cdot s_1 \ t \cdot s_2 \ \dots]$ is a row vector whose elements are dummies that signal if the subject status has changed in year 2007 compared to 2003. For example, imagine a subject becoming a parent in the interval 2003–2007. The corresponding value of his/her dummy status, *s*, is equal to 1. The interacting term $t \cdot s$ (included into the vector $t \cdot \mathbf{s}'$) takes value zero in 2003 (*t* = 0) and value 1 in 2007 (*t* = 1). Thus, the scalar product between the vector $t \cdot \mathbf{s}'$ and the vector of coefficients $\boldsymbol{\delta}$, given by $t \cdot \mathbf{s}' \boldsymbol{\delta}$, measures the interaction effect of status and time.⁵

The vectors \mathbf{x}_i and \mathbf{z}_i are, respectively, time-invariant and time-varying control variables, with corresponding parameters $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$. Control variables, referring to the

⁵In our model, we consider several changes of status referred to specific life course transitions, as in Baxter et al., 2008 (see also [15], pp. 146–151). Conley and Taber [5] found that standard methods generally used to perform inference in *DID* models with more transitions, or with the same transition observed in different groups, are not completely appropriate, and lead to underestimation of the standard errors size. Nevertheless, we decided to estimate more transition effects in our analysis in order to obtain a richer model specification.

subject or to the couple,⁶ are given by education, the region of residence, the place of residence, age, years of work experience, the hourly wage, the economic value of the house of residence (as a proxy of household wealth). Other covariates are given by gender attitudes measures and dummies representing life course events transitions (Table 2). We adopt the logarithms of weekly working hours as dependent variables, while explanatory variables are included in the regression using original values.⁷ In doing so, coefficients of these regressors can be considered as “semi-elasticities” measuring the impact on the dependent variable in percentage terms. This allows us to better evaluate the joint influence of two or more regressors by adding together the respective estimated coefficients.⁸

For each equation, we assume that the error terms, u , are distributed with zero mean, and that the covariances between the error terms of each equation (cross-sectional endogeneity) and across time (correlation between repeated observations) are non-zero. Therefore, the covariance matrix, Ω (similar to the one used in a *SURE* model), of the error terms of the four equations may be specified as follows:

$$\Omega = \overset{8 \times 8}{\Sigma} \otimes \overset{n \times n}{\mathbf{I}} = \begin{bmatrix} \overset{4 \times 4}{\Sigma_{00}} \otimes \overset{n \times n}{\mathbf{I}} & \overset{4 \times 4}{\Sigma_{01}} \otimes \overset{n \times n}{\mathbf{I}} \\ \overset{4 \times 4}{\Sigma_{10}} \otimes \overset{n \times n}{\mathbf{I}} & \overset{4 \times 4}{\Sigma_{11}} \otimes \overset{n \times n}{\mathbf{I}} \end{bmatrix} \quad (5)$$

where \mathbf{I} is the identity matrix and Σ is a block-matrix, in which Σ_{01} and Σ_{10} are sub-matrices of dimension 4×4 showing in the diagonal n constant elements given by the covariances between the error terms across time referring to the same subject in each equation (across time correlation). Σ_{00} and Σ_{11} are symmetrical sub-matrices including variances of error terms of each equation in the diagonal and covariances between the equations elsewhere (correlation “between” individuals).

The residuals of an *OLS* regression of Eqs. (1)–(4) are used to estimate the error covariance matrix as previously specified. The residual-based estimation of the error covariance matrix then allows us to run iteratively a *GLS* procedure. In this way, the bias in estimates due to omitted variables can be corrected. The estimation procedure here adopted is implemented by combining ad hoc different STATA 11 commands.⁹

However, a further estimation problem occurs because 47 women who were employed in the year 2003 interrupted their working activity in the period between 2003 and 2007. Therefore, information on women’s paid work hours were censored

⁶To simplify the model, variables referred to the partner are not included as regressors.

⁷One exception is the logarithm of hourly wage, whose coefficient measures the wage elasticity of labour supply.

⁸Note however that the logarithmic transformation of some variables (monotonic transformation, in any case) may introduce small distortions.

⁹An iterative *GLS* procedure has been implemented by writing a “do” file in STATA. In particular, we employ, at each iteration, the residual-based estimate of error terms covariances to correct linear regressions.

in 2007 for these 47 women. As a consequence, a selection bias may affect paid work estimation. The solution here adopted to avoid the selection bias is to run a *Tobit Random-Effects* [15] regression only for women's paid work hours (Eq. 1) at the first step of the *GLS* iterative procedure. This approach is suggested by Lee [11] as a two-step estimator of a simultaneous equations model with limited and censored dependent variables.¹⁰ The assumptions underlying this model are that the error terms are normally distributed with zero means and a covariance matrix that is constant over all observations. In particular we use a *Tobit Random-Effects regression* to take into account the panel structure of dataset.

To simplify the model and to reduce the number of regressors, we estimate the effects of parenthood transition by using three distinct sample stratifications corresponding to different birth orders. First, we take into account the couples who experienced the transition to the first birth, using couples without children in 2003 as reference group. Second, we estimate the effect of transition to a second-order birth, considering couples with one child in 2003 as a reference group. Third, we consider couples who experienced the transition to parenthood (whatever the order) in the period 2003–2007.¹¹

4 Empirical Analysis and Estimation Results

In order to verify empirically how omitted factors may affect model estimation for cross-sectional endogeneity, we estimate preliminarily the correlation matrix of *OLS* residuals “between” equations. We find, for women, that the correlation coefficient between domestic work and the paid work *OLS* residuals is equal to -0.12 in the year 2003 and -0.31 in the year 2007; while the correlation between residuals of domestic work of women and paid work of men is equal to 0.095 in the year 2003 (0.079 in the year 2007).

Subsequently, we stratify the sample for different orders of birth and we estimate the model adopting our *SURE-DID* procedure. In Table 1, we report only coefficients referring to the transition to parenthood and log-hourly wage. In particular, parenthood coefficients refer, respectively, to all orders of birth, first birth and second birth. In Table 2, the estimated coefficients of the model considering all orders of birth are presented extensively¹² (coefficients measuring the effect of a birth are highlighted). The effect of the birth of a child between 2003 and 2007 is measured by the estimated *DID* coefficients (t^*Birth), that may be interpreted as a percentage variation of the dependent variable due to the transition to parenthood.

¹⁰Note that in several studies *SURE* models with censored dependent variables are estimated using the Lee method as, for instance, the demand systems models [8, i.a.].

¹¹The couples who had two (or more) children between 2003 and 2007 (12 couples in total) are included only in the sample containing all orders of births.

¹²For sake of brevity, here we do not present extensively the estimates of the stratified analysis by order of birth, referring, respectively, to the first-order birth and to the second-order birth transition.

Table 1 Estimates of the *SURE-DID* model considering the sample with all orders of birth and sub-samples for transition to first and second birth

Dependent variable (log):	Women		Men	
	Paid work coeff.	Unpaid work coeff.	Paid work coeff.	Unpaid work coeff.
<i>Transition to parenthood (all orders, no. 562 couples; no. 1124 subjects; no. 122 transitions)</i>				
<i>t</i> (Dummy: 0 = 2003; 1 = 2007)	-0.47***	0.05	0.01	0.15*
<i>W</i> (log real-hourly wage)	1.36***	-0.83***	0.13	-0.68*
<i>Birth</i>	0.14*	-0.31***	-0.03	0.09
<i>t*Birth (DID coeff.)</i>	-0.40***	0.55***	-0.02	0.05
<i>Transition to the first child (no. 86 couples; no. 172 subjects; no. 45 transitions)</i>				
<i>t</i> (Dummy: 0 = 2003; 1 = 2007)	-0.22	0.10	0.02	0.27
<i>W</i> (log real-hourly wage)	0.55	-1.78***	-0.11	-1.08
<i>Birth</i>	-0.01	-0.21	-0.05	0.15
<i>t*Birth (DID coeff.)</i>	-0.36	0.55*	0.03	0.22
<i>Transition to the second child (no. 182 couples; no. 364 subjects; no. 65 transitions)</i>				
<i>t</i> (Dummy: 0 = 2003; 1 = 2007)	-0.57***	-0.04	-0.04	0.30*
<i>W</i> (log real-hourly wage)	1.25**	0.23	0.33	-0.49
<i>Birth</i>	0.19	-0.23*	-0.02	0.30*
<i>t*Birth (DID coeff.)</i>	-0.51**	0.57***	0.04	-0.11

P-value: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$

Only coefficients concerning to parenthood are here shown. In particular, the difference-in-differences coefficients are written in bold

As shown in Tables 1 and 2, the estimated *DID* coefficient reveals that the impact of transition to parenthood (all orders) on the labour hours of Italian married women is negative and equal, in percentage, to -40% . We also obtained a strong positive influence of wage elasticity, measured by the coefficient of log (real-hourly) wage, in the female paid work equation. However, this income effect cannot compensate for the negative effects of the birth of the child.¹³ The negative impact of motherhood on paid work is even stronger (-36%) if we consider women who experienced the transition to a second birth (Table 1).

Estimation results in Table 1 also show that women's hours of unpaid activity increase by 55% as a consequence of transition to parenthood. Considering birth order, women's unpaid activity still increases by 55% after the first child and by 57% after the second child. Note that for men, contrary to what happens for women, market and domestic working activities seem to be less sensitive to the transition to fatherhood.

Analyzing the effect of gender attitudes proxies we found that several respondents changed their opinion between 2003 and 2007 regarding the statement that, if spouses get a divorce, child custody should be assigned to the mother (see dummy:

¹³To compensate for the reduction of paid work caused by the transition to motherhood ($-0.40 + 0.14 = -26\%$) and given the wage elasticity coefficient equal to 1.36 (Table 1), women's hourly wage should increase by a percentage of 19% (obtained by the ratio $0.26/1.36$).

Table 2 Estimates of the *SURE-DID* model considering the sample with all orders of birth and sub-samples for transition to first and second birth

Dependent variable (log): Work	Women		Men	
	Paid coeff.	Unpaid coeff.	Paid coeff.	Unpaid coeff.
Constant	1.91*	3.27***	4.43***	-1.57
Controls				
<i>Expe</i> = work experience (TV)	0.03***	-0.01**	0.01	-0.01
<i>Age</i> = age (TV)	-0.08*	-0.05	-0.05*	0.11*
<i>Age</i> ² = age ² (TV)	0.00	0.00	0.00	0.00*
<i>Edu</i> = education level—years of schooling (TI)	0.03***	-0.01	0.00	0.01
<i>Urban</i> = Dummy metropolitan area Yes = 1 (TI)	-0.02	-0.11**	0.01	0.02
<i>House</i> = House economic evaluation indicator	-0.02	-0.06**	0.02	-0.03
<i>Area</i> = Dummy: South = 1 (TI)	-0.01	0.10*	0.00	-0.21***
<i>t</i> = (Dummy: 0 = 2003; 1 = 2007)	-0.47***	0.05	0.01	0.15*
<i>W</i> = log (real)hourly wage (TV)	1.36***	-0.83***	0.13	-0.68*
<i>Mean_W</i> = Mean of log (real) hourly wage (TI)	-0.28	1.28***	0.08	1.05*
Gender Attitude				
<i>Child custody</i> = Dummy: child custody if spouses divorce 1 = Mother (TV)	0.02	0.05	0.02	-0.18*
<i>t*Child custody</i>	0.09	-0.16*	0.08	0.13
<i>Eldercare</i> = Dummy: elder-care duty in the family 1 = female (TI)	-0.01	-0.01	-0.07*	-0.13
<i>Partneraid</i> = Dummy: If high order birth depends on partner commitment: 1 = Yes (TI)	0.05	0.04	0.03	0.24***
Birth Transition (transition to parenthood, all orders)				
<i>Birth</i>	0.14*	-0.31***	-0.03	0.09
<i>t*Birth</i>	-0.40***	0.55***	-0.02	0.05
Work Transition				
<i>Full_part</i> = Dummy: 1 if work Full Time 2003 and Part Time 2007	0.13	-0.09	0.09*	-0.26**
<i>Part_full</i> = Dummy: 1 if work Part Time 2003 and Full Time 2007	-0.19	-0.05	-0.02	-0.46***
<i>t*Full_part</i>	-0.07	0.16	-0.14*	0.18
<i>t*Part_full</i>	0.62***	-0.01	0.01	0.56**

(continued)

Table 2 (continued)

Dependent variable (log): Work	Women		Men	
	Paid coeff.	Unpaid coeff.	Paid coeff.	Unpaid coeff.
Union Status Transition				
<i>Dummy Cohabiting2003 and 2007 = 1 (yes)</i>	-0.01	-0.01	-0.25***	0.43*
<i>Dummy Cohabiting2003-married 2007 = 1 (yes)</i>	0.00	-0.31	-0.07	0.63
<i>t* Cohabiting2003 and 2007</i>	0.32	-0.18	0.23*	-0.58*
<i>t* Cohabiting2003 and married 2007</i>	0.33	0.02	0.10	-0.43
Help Transition (Received)				
<i>Help elders = Dummy: 1 if the family received help caring for the elderly</i>	0.23	1.11*	0.25	-1.25
<i>Help children = Dummy: 1 if family received help for childcare</i>	-0.01	0.17***	0.04	0.23**
<i>Help domestic = Dummy: 1 if family received help for domestic chores</i>	-0.15	0.08	-0.01	0.16
<i>t* Help elders</i>	0.07	-1.45**	-0.10	0.76
<i>t* Help children</i>	0.01	-0.23**	-0.06	-0.14
<i>t* Help domestic</i>	0.27*	0.07	-0.05	-0.18
R ²	0.23	0.10	0.05	0.10

TI: Time invariant, *TV*: Time variant. The difference-in-differences coefficients are written in bold
P-value: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$

Child custody in Table 2). More specifically, 61 women and 20 men changed their opinion from 1 = agreement to 0 = disagreement. For this reason, we model the latter dummy taking into account the change over time by introducing a specific *DID* parameter (dummy: t^* *Child custody*). Females changing their opinions seem to change their allocation of time, working less at home and more in the labour market.

Estimation results of Table 2 show also that help received (informal or paid) by the family in care activity and domestic work generally contribute to a reduction in a woman's domestic work. Namely, estimation results reveal a strong reduction in domestic work as an effect of help received for caring for the elderly and children.

5 Conclusion

Analyzing the relationship between transition to parenthood and the partners' allocation of time between paid and unpaid activities, we found that the Italian women's earnings ability effect on their paid work supply (measured, in Table 2, by

education, experience and wage coefficients) is positive and significant. However, the negative impact of the birth of a child is stronger and leads the partners to a more gendered division of market and domestic work in the family. Moreover, as shown in other recent studies [4], we found that changes in the couple's fertility, working status, marriage and changes in help received in care and domestic activities do affect woman's time, while man's time remains mostly unaffected.

In this study we adopt a *SURE-DID* procedure. The *DID* specification permits us to disentangle in parenthood transition (a) the actual life event effect from (b) the selective nature of those who experience the life course transition. The *SURE* stochastic specification of the error terms here adopted lets us check for endogeneity effects of partners' behaviour and over-time correlation between repeated observations. Compared to fixed effect or change-score models adopted in other longitudinal studies [9, 10], our *SURE* procedure lets us obtain estimations also for the effect of time-invariant variables. In addition, our methodological approach allows us to estimate simultaneously for both partners the impact of life course events (e.g. the birth of a child) on the allocation of time in paid and unpaid activities.

Some shortcomings of our study concern the limited number of couples (even if this is quite a common characteristic of socio-demographic longitudinal studies) and the selective nature of our sample, as only couples working in 2003 and still together in 2007 are considered. Finally, some caution is needed, as for most social studies, for the notion of causality. Thus, expressions such as "influence", "effect" or "depend on" should be interpreted accordingly.

Acknowledgments We would like to thank participants at the 46th Scientific Meeting of the Italian Statistical Society—Session "Demographic methods and models"—for their useful comments.

References

1. Aassve, A., Burgess, S., Propper, C., Dickson, M.: Employment, family union and childbearing decisions in Great Britain. *J. R. Stat. Soc. A Sta.* **169**(4), 781–804 (2006)
2. Abadie, A.: Semiparametric difference-in-differences estimators. *Rev. Econ. Stud.* **72**(1), 1–19 (2005)
3. Baxter, J., Hewitt, B., Haynes, M.: Life course transitions and housework: marriage, parenthood, and time on housework. *J. Marriage Fam.* **70**, 259–272 (2008)
4. Campolo, M.G., Di Pino, A.: An empirical analysis of women's working time, and an estimation of female labour supply in Italy. *Statistica* **72**(2), 173–193 (2012)
5. Conley, T.G., Taber, C.R.: Inference with "Difference in Differences" with a small number of policy changes. *Rev. Econ. Stat.* **93**(1), 113–125 (2011). doi:[10.1162/REST_a_00049](https://doi.org/10.1162/REST_a_00049)
6. De Santis, G.: The monetary cost of children. Theory and empirical estimates for Italy. *GENUS* **61**, 161–183 (2004)
7. Grilli, L., Rampichini, C.: The role of cluster sample means in multilevel models: a view of endogeneity and measurement error issues. *Methodology-Eur.* **7**(4), 121–133 (2011). doi:[10.1027/1614-2241/a000030](https://doi.org/10.1027/1614-2241/a000030)
8. Heien, D., Wessels, C.R.: Demand systems estimation with microdata: a censored regression approach. *J. Bus. Econ. Stat.* **8**(3), 365–371 (1990)

9. Johnson, D.R.: Alternative methods for the quantitative analysis of panel data in family research: pooled time-series models. *J. Marriage Fam.* **57**(4), 1065–1077 (1995)
10. Johnson, D.R.: Two-wave panel analysis: comparing statistical methods for studying the effects of transitions. *J. Marriage Fam.* **67**(4), 1061–1075 (2005)
11. Lee, L.F.: Simultaneous equations models with discrete and censored dependent variables. In: Manski, P., McFadden, D. (eds.) *Structural Analysis of Discrete Data with Econometric Applications*, pp. 346–364. MIT Press, Cambridge (1978)
12. Mills, M., Mencarini, L., Tanturri, M.L., Begall, K.: Gender equity and fertility intentions in Italy and the Netherlands. *Demogr. Res.* **18**(1), 1–26 (2008). doi:[10.4054/DemRes.2008.18.1](https://doi.org/10.4054/DemRes.2008.18.1)
13. Sanchez, L., Thomson, E.: Becoming mothers and fathers: Parenthood, gender and the division of labor. *Gender Soc.* **11**, 747–772 (1997)
14. Srivastava, K.V., Giles, D.: *Seemingly Unrelated Regression Equations Models*. Marcel Dekker, Basel/New York (1987)
15. Wooldridge, J.M.: *Econometric Analysis of Cross Section and Panel Data*, 2nd edn. MIT Press, Cambridge (2010)

Do Rational Choices Guide Family Formation and Dissolution in Italy?

Gustavo De Santis and Silvana Salvini

Abstract

As social pressure to adhere to accepted standards recedes, individuals are freer to choose their preferred family arrangement. But their choices do not always follow a strictly rational approach: rather, they appear to be guided by a trial-and-error logic, which implies contradictions, loss of efficiency and, occasionally, low personal satisfaction. Modern welfare states, while supporting freedom of choices (which includes the possibility of change: e.g. divorce, or medically assisted fertility in one's late years), must combine this with other targets, ranging from the empowerment of women to the protection of the weak (children, especially), to a system of incentives that eventually ensures socially acceptable outcomes, including a sufficient level of fertility. The now undisputed primacy of the individual will not destroy families, but it will deeply transform them, and increase their heterogeneity: relationships will increase in number but decrease in duration and intensity.

1 The "Pursuit of Happiness"

"The Pursuit of Happiness" (yes, misspelt) is a 2006 movie directed by Gabriele Muccino, based on the autobiography of Chris Gardner, once a homeless, now a successful stock broker. It is not totally clear whether the happiness that Gardner eventually achieves depends on his success in economic or in family matters. Both, probably, since both began badly and ended well. Chris, who never knew his father

G. De Santis (✉) • S. Salvini

DiSIA Dip. di Statistica, Informatica, Applicazioni "G. Parenti", Viale Morgagni 59, 50134 Florence, Italy

e-mail: gustavo.desantis@unifi.it

(well, yes, much later: at the age of 28), spent his childhood in poverty with his mother, his older sister, two more stepsiblings and an often drunk and violent stepfather. Chris's mother was imprisoned twice: first, when his stepfather falsely accused her of cheating on welfare, and then when she took revenge and tried to kill him by burning their house down—with him inside. Chris and his siblings went to foster care, but the beloved uncle they lived with drowned in the Mississippi River when Chris was just 9 years old. As an adult, Chris got married, but then started an affair with a younger girl, who got pregnant and gave birth to his son, Chris Jr. Understandably, his wife divorced him, but his new partner too abandoned him, shortly after, when Chris was imprisoned for debts. Fresh out of jail, he and his son, just a toddler then, lived as homeless for almost a year, with no money and no family support. Scarce prospects indeed for happiness, and yet . . .

The pursuit of Happiness (but this once spelt correctly) is also listed, after Life and Liberty, among the unalienable rights of all men, according to the Declaration of Independence of the United States. What Thomas Jefferson, composer of the original draft, back in 1776, intended by "Happiness" is not totally clear, but he probably had in mind what Inglehart [1] would later call "materialist" values: basically economic and physical security. In those hard times, these goals required strong associations between individuals, based on hierarchy, and on a set of structures (village, church, . . .), among which was the family, with its undisputed head, and its formal and rigid rules. Individual values (liberty, autonomy, independence, . . .) did not matter much, then.

Much later, however, and this holds especially for the cohorts born after World War II, economic security led to the emergence, and eventually the predominance, of "post-materialist" values (Fig. 1). This profound change of priorities did not take place abruptly: it followed the renewal of generations because "people are most likely to adopt those values that are consistent with what they have experienced first-hand during their formative years" [2, p. 132], and those who grew up "with the feeling that survival can be taken for granted" attach more importance to such things as autonomy, self-expression, quality of life, freedom of speech, relevance of original ideas, giving people more say in important government decisions, etc". It is an "aspect of a [. . .] broader process of cultural change that is reshaping the political outlook, religious orientations, gender roles, and sexual mores of advanced industrial society. The emerging orientations place less emphasis on traditional cultural norms" [2, p. 138], and, for what concerns us here, lead to "gender equality . . . tolerance of outgroups, including foreigners, gays and lesbians [while] younger cohorts become increasingly permissive in their attitudes toward abortion, divorce, extramarital affairs" [2, pp. 139–140]—with the exception of Gardner's wife, of course.

Note that this generation-based approach of shifting preferences is not far from Max Planck's earlier [3, p. 22] and disenchanting opinion that "a new scientific truth does not triumph by convincing its opponents and making them see the light, but rather because its opponents eventually die, and a new generation grows up that is familiar with it". Indeed, Fig. 1 suggests that values and orientations (including

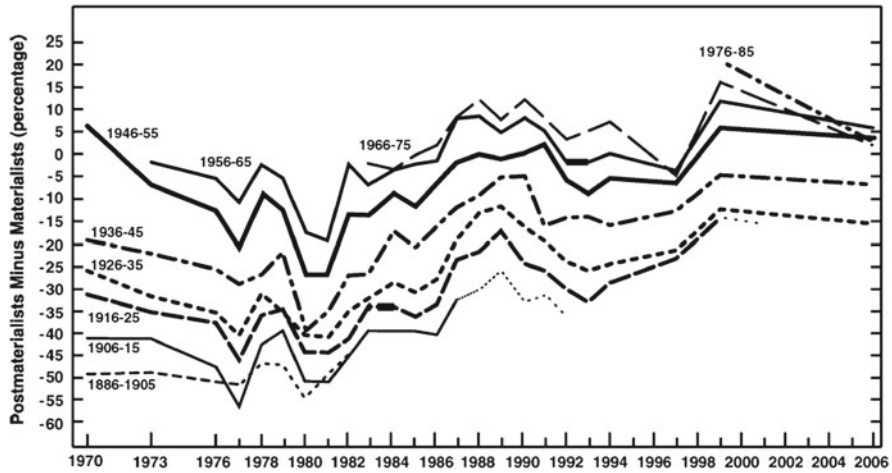


Fig. 1 Cohort analysis: % post-materialists minus % materialists in six west European societies, 1970–2006. *Source:* Based on combined weighted sample of Eurobarometer surveys and World Values Surveys in West Germany, France Britain, Italy, the Netherlands and Belgium, in given years, using the four-item materialist/post-materialist values [2, p. 135]

prejudices) do not vary much with age, but do evolve with generations—in the case of Fig. 1, towards post-materialist values.

This is the “Silent Revolution”, in Inglehart’s words, which echoes, both in terminology and substance, the “Quiet Revolution” that had taken place in Quebec (CAN), a few years before, in the 1960s. Its most salient features are rapid secularization and the creation of the modern welfare state, whereby health care and education, previously in the hands of the Roman Catholic Church, passed under the responsibility of the provincial government. Here, too, survival could finally be taken for granted, (religious) authority started to dwindle, and new perspectives opened up for individual preferences, in various domains, including family matters.

2 Between Tradition and Rationality

Under the impulse of Gary Becker [4], the *New Household Economics* (NHE) appeared on the scene, and started to apply the paradigms of economic logic to demography: fertility and the demand for children (where “quantity” is apparently being traded off for “quality”), altruism in the family, intra-household bargaining and gendered division of labour, partnering and re-partnering, and similar topics have been studied under a different light, since then. Basically, the idea now is that all demographic choices are taken after careful examination of the available information, given constraints (in terms of time and money) and personal orientation, with the aim of maximizing some individual utility function. This, to be sure, still leaves a number of important issues open: for instance, how to reconcile the not

necessarily converging preferences of various family members, which is definitely more difficult if all have the same rights and are no longer forced to follow the lead of the “household head”, or how to make sure that rational individual choices do not lead to undesirable aggregate outcomes, typically because of externalities, insufficient information or lack a proper system of incentives and disincentives. Fertility is a good example here: individuals may prefer to have too few children (as in most developed countries) or too many (as in sub-Saharan Africa), and the question then arises of how to induce them to have (roughly) the number that is “right” for society—if this ideal number exists at all.

What does the *SDT* (*Second Demographic Transition*) theory add to all this? Not much, in our opinion—but the vast echo that it had, and still has, in the specialized literature reveals that demographers worldwide think differently. It shares with the *NHE* the idea that choices are taken by individuals, and not by institutions (state, church, family, tradition, etc.). But it insists on the importance of the “cultural shift”: once freed from “authority”, individuals do not simply pursue traditional goals with a new rationality; rather, they start to develop personal goals [5], and this explains the increasing variety of (individual) demographic choices, family forms and living arrangements [6, p. 137].

More generally, however, the contrast between the *NHE* and the *SDT* theory implicitly refers to a possible opposition between “economics” and “culture”. What is it that (mainly) drives our demographic decisions: is it our personal advantage, that we coldly calculate, or is it the socio-cultural environment, that in part shapes our preferences, and in part forces us along an “expected”, normative path, from which we would like, but do not dare, to deviate for fear of social disapproval? This may apply to *de facto* unions [7], out-of-wedlock births, frequency of contacts with one’s elderly parents [8], etc. In part, this may be a false problem: internalization of norms and formation of desires and goals come first, and this is what the *SDT* theory is about. Once these targets are defined, the problem arises of how to best achieve them, and this is what the *NHE* deals with. However, with regard to family processes, there are a few assumptions that the *NHE* and the *SDT* theory have in common, that are relevant for family formation and dissolution, and that it might be worthwhile to reconsider, at least in part, because they do not seem to be fully consistent with empirical evidence.

3 What Do We Mean by Rational?

“Dynamic economic models are based on the forward looking behavior of economic agents. In the context of life-cycle models, an individual’s consumption and savings decision depends on her subjective beliefs about future interest rates, wage rates and the likelihood of dying. According to these models, individuals have beliefs about such variables and use these beliefs to make decisions today. Until recently common practice in such studies was to assume rational expectations implying that the individuals’ beliefs are given as objective probability distributions. The use of objective distributions is by now put into question by numerous researchers who

suggest to directly measure subjective expectations and to evaluate the consequences of deviations of subjective expectations from their objective counterparts” [9].

In other words, people may take their decisions on the basis of false premises. In part, this happens because they do not have enough information, or do not know how to correctly interpret what they know. For instance, when asked about their own future survival prospects, on average “people under-estimate how long they are likely to live by over 5 years. They tend to ignore expected mortality improvements” [10, p. 31]. At the same time, “people are optimistic: they think they will live longer, on average, than people of their own age and sex: by about by 1.19 years (males) and 0.76 years (females).” [10, p. 32]. This may have huge implications on pension policies, for instance: the idea that people know what is best for them, e.g., in terms of age at retirement, is considered almost a tautology, nowadays. But closer inspection reveals that “policy makers . . . cannot assume that people share a rationale to prepare for a retirement of a realistic length” [11, p. 198].

As for couple formation, how else can we interpret the fact the pre-marital cohabitation reduces [12] or, at the very best, does not increase [13] the solidity and the average length of marriages? Partners who have had the opportunity to test each other do not seem to benefit from this experience: why? A possible answer is that interpreting the available information (on the labour market, on the partner, and on a number of other topics) proves too difficult for many, maybe even for the majority of people.

Secondly, people are rarely capable of forecasting the future. If this holds for professional demographers, when they try to anticipate the likely course of populations (i.e., when they do what they should in principle be best equipped to), what else should we expect of common people faced with experiences that are in most cases totally new to them? So, when asked about their pension plans, in a sample of mature Dutch workers, those who think that they will live longer than average also (correctly) state that they intend to retire later. In practice, however, they do not: their age at retirement is just average [14]. The same survey reveals that all the interviewed grossly overestimate their age at retirement, by about 1.6 years, despite the fact that they are all mature workers, aged between 50 and 64, and therefore relatively close to retirement and well aware of the labour market situation, both in general and with reference to their specific case.

It should not come as a surprise, then, that an even greater incapability can be discerned in several other demographic domains, too, closer to our topic of interest. Margolis and Myrskylä [15, p. 48], for instance, conclude that “people seem to poorly predict how children affect their lifestyle and underestimate the costs of children”. On the other hand, assisted reproductive technologies (e.g. artificial insemination, or in-vitro fertilization) are on a spectacular rise, in France (where about 5 % of all births, at the beginning of the twenty-first century, fall in this category) and elsewhere [16]: this happens in part because people are now less prone to accept whatever “God sends” (religious faith is dwindling, remember?), which includes childlessness; in part, because fertility has been delayed so much (age at first birth is now well over 31 years in Italy; [17]) that some women later find it difficult to achieve their desired, even if relatively low, fertility. When it comes to

the realization of fertility intentions, less than 50 % of the interviewed, in a Dutch panel, end up exactly with the number of children they had intended to have when they were 26 years old, and this holds especially for women [18].

But incapability of anticipating the future may be only part of the problem. The other part—and this is the third point that we want to make—derives from the fact that people are not consistent in their aspirations, and not on trivial matters only. With regard to childbearing intentions, for instance, Iacovou and Tavares [19, p. 117] agree with Liefbroer that “many people simply change their mind”, typically wanting fewer children later than they did initially, especially after having experienced parenthood, which frequently turns out to be less rewarding than originally imagined (see “incapability of anticipating the future”). Similarly, Testa [20, p. 5] warns that “reproductive ideals and intentions, as well as actual fertility, are developmental by nature and change over the individual’s life course . . . the adjustments of fertility goals over the life course tend to occur mainly downward in response to different factors and events, one of which is of particular importance, i.e., the transition to a first or a higher birth order child”.

The rise in demand for adopted children that was observed up to a few years ago [21], and that is now being progressively substituted by a demand for medically assisted reproduction (in those countries where this is permitted by law, technology and resources; see again [16]) falls in part in this category, because, up to a certain point at least, it comes from people who had formerly decided not to have children (or had light-heartedly accepted this possibility), but later change their mind.

All the surveys on happiness and life satisfaction consistently show that those who are married are the happiest; then come those who cohabit, then those who are single, and last those who used to live in couple, but are now alone, because of widowhood or marriage breakdown [15,22–25]. Yet, the share of single-person households is everywhere on the rise: in Italy, for instance, up from 9 % in 1901 to almost 30 % now. And it would be even higher, had it not been for immigrants, whose households are typically larger than ours. In part, this depends on a few structural factors, and notably population ageing, but delays in couple formation and increases in divorce and separations, too, play a prominent part. Separations, for instance, are high (almost 30 % of marriages end in a separation, in Italy; between 40 % and 50 % in most of Europe) and on the rise; and couple dissolution is even higher among the cohabiting. Besides, the divorced and the separated, especially when they have children (which happens frequently), also face economic difficulties and time constraints: this increases their stress [26] and is a non-secondary cause of their low levels of happiness. It is not totally clear whether marriage breakdown depends on incapability to interpret the initially available information (on the partner), inability to predict the future (changes in the partner’s traits), or inconsistency (change of preferences). But it is invariably associated with a feeling of failure and loss of confidence (in oneself, in the future, etc.: [23]), which, we contend, is amplified by the fact that the idea that we all proceed by trial and error (error, especially) is not widely accepted, yet. Actually, to the best of our knowledge, it does not even seem to have ever been discussed in the demographic literature.

4 Making and Unmaking Families by “Trial and Error”

If people make mistakes and change their mind, there are all reasons to expect a growing variety of family forms, and strong individual mobility among them—unless, of course, something prevents it: for instance, law provisions, tight social control or lack of economic resources. Indeed, the very notion of marriage has changed over time: the label may be the same, but what was once “until *death* do us part” is now “until *we* do us part”. Where is the difference, then, from a simple cohabitation? Everywhere illegitimate children are no longer discriminated (the very term, illegitimate, is now abolished or anyway avoided), and unmarried partners are nowadays entitled to a few (and growing) rights. Breaking a marriage is, of course, more demanding, but the tendency seems to be towards quicker and cheaper solutions, in several (and increasing) cases regulated by prenuptial agreements, where explicit provisions are introduced to regulate the possible termination of marriage.

It is not unlikely that families will eventually become a non-permanent attribute of individuals, with relatively frequent changes in arrangements, partners, houses, etc. Should we worry? Let us rapidly review the potential shortcomings of this evolution. One is that fertility may be (further) depressed by this system of “light” partnering. This is possible, of course, but in the past 25 years fertility has been much lower in southern Europe, with its “strong” families, than in northern Europe, where the system of “modern” families first emerged. Fertility is sustained much more by female empowerment (on the labour market, especially) and active family support (child care facilities, notably) than by prohibition to divorce, or any other constraint.

Children of broken families suffer, and it not easy to guarantee a fair distribution of custodial rights and obligations after couple breakdown. True, but once partners can no longer stand each other, forcing them to remain together may not be the best solution, not even for the children [27]. And if marriage termination is relatively painless for parents (e.g. thanks to prenuptial agreements), why should partners quarrel? This should also ease subsequent cooperation, instead of conflict, in rearing children. Besides, as Leridon [28] suggests, fertility, too, is slowly evolving towards individualization: even without clones, an efficient system of donors (of both sperm and oocytes) will permit individuals to have the children they want, even without partners. If this is indeed the future that awaits us, children of broken couples should not constitute a major problem, especially once they cease to be, and to be considered, an exception. And if Chris Gardner made it, out of *his* family, why shouldn’t anybody else?

Another possible objection is that families have traditionally been responsible for caring for weak members, e.g. the old and the sick. Who will, if and when families lose their strength? Good question, but let us not forget that this “nurse service” has thus far been provided essentially by women, precisely because of their “lesser” role in society, and in the labour market, and this is no longer going to be the case. It simply means that care costs, thus far largely hidden (and imposed on women) will

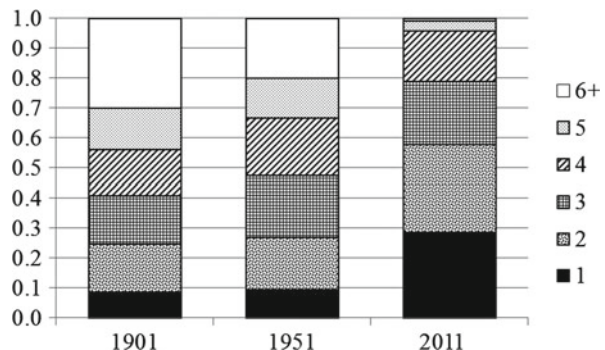
emerge in full, and will have to be faced, somehow. Not an easy or cheap passage, surely, but so was the abolition of slavery: who would object to it, nowadays?

Finally, and more generally, living in large families generates economies of scale. A future of individuals forming small families, and possibly only for short periods will prove costly. True: but this is already happening. Just as an illustrative example, consider that between 1901 and 2010 the population of Italy almost doubled: the index number is 192 (if 1901 = 100). In terms of equivalence scale, however, costs have increased considerably more: to somewhere between 231, with Carbonaro's scale [29], and 264, with OECD's [30] square root scale, precisely because the average family size has shrunk (Figs. 2 and 3). Per capita income has in the meantime increased much more (the index number is close to 2000: see <http://www.ggd.net/MADDISON/oriindex.htm>), which simply means that we have decided to use part of our extra riches to buy something that we value (privacy and intimacy), even if it is expensive.

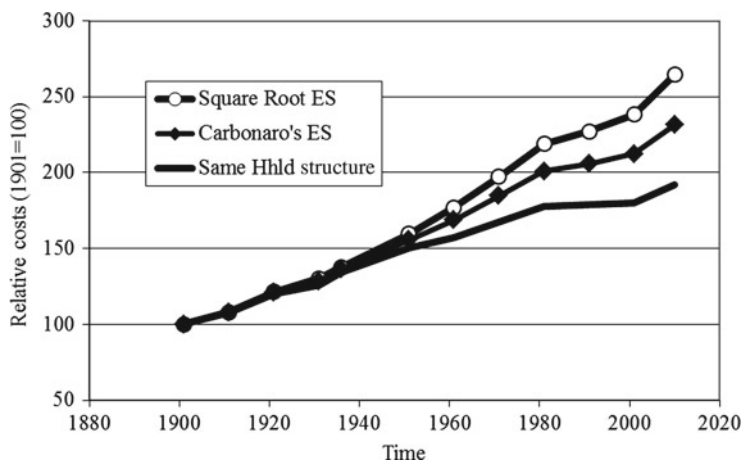
Together with intimacy (see also lines 1–5 of Table 1), we have also bought variety: single parents, unmarried couples and reconstituted families, for instance, are on the rise, although most of these changes can be documented only for the past 20 years or less. There are also indications that constraints, too, and not only choices, play their part: the number of young adults, aged 18–30, who still live with their parents, and are therefore classified as “children” in population surveys (Table 1), and whose growing share depends in large part on the unfavourable labour market conditions of this age group [32].

The family is changing in Italy, and probably also elsewhere, and it has evolved from an institution to an individual attribute. Within the limits imposed by the respect of other's needs, and by personal resources, and unless major structural changes occur, this “Quiet Revolution”, which has only just begun, will most likely continue. Individuals will thus be left to themselves in their “pursuit of happiness”, but, in the process, mistakes, changes of opinion and direction, disappointments and costs, of all kind, will be the rule rather than the exception.

Fig. 2 Distribution of households by size (Italy, selected years)



Source: Population censuses



Source: Own elaboration on census data

Fig. 3 Estimated costs of larger population and shrinking household size (Italy, 1901–2011)

Table 1 Selected characteristics of households and families in Italy (1994–2007)

	<i>Out of 100 . . .</i>	1994–1995	1998–1999	2002–2003	2006–2007
Singles	<i>Households</i>	21.1	22.2	25.3	26.4
Households with 5+ members	<i>Households</i>	8.4	7.7	6.8	6.2
Extended households	<i>Households</i>	5.1	5.5	5.3	4.8
Couples with children	<i>Nuclear families</i>	62.4	60.8	58.9	56.8
Couples without children	<i>Nuclear families</i>	26.7	28.1	29.2	30.6
Single parents	<i>Nuclear families</i>	10.9	11.1	11.9	12.7
Unmarried couples	<i>Two-sex couples</i>	1.8	2.4	3.9	4.6
Reconstituted families	<i>Two-sex couples</i>	4.1	3.9	4.8	5.6
Unmarried children, 18–30 years	<i>People aged 18–30</i>	69.5	72.4	72.7	72.8

Source: Istat [31]

References

1. Inglehart, R.: The silent revolution in Europe: intergenerational change in post-industrial societies. *Am. Polit. Sci. Rev.* 65(4), 991–1017 (1971)
2. Inglehart, R.: Changing values among western publics from 1970 to 2006. *West Eur. Polit.* 31(1–2), 130–146 (2008)
3. Planck, M.: *Scientific Autobiography and Other Papers*. New York 1949 (original edition 1948). As cited in T.S. Kuhn, *The Structure of Scientific Revolutions*)
4. Becker, G.S.: *A Treatise on the Family*. Harvard University Press, Cambridge, MA (1981)
5. Sobotka, T.: The diverse faces of the second demographic transition in Europe. *Demogr. Res.* 19(8), 171–224 (2008)

6. De Rose, A., Vignoli, D.: Families “all’italiana”: 150 years of history. *Riv. It. Di Econ. Demogr. E Stat.* **65**(2), 121–144 (2011)
7. Di Giulio, P., Rosina, A.: Intergenerational family ties and the diffusion of cohabitation in Italy. *Demogr. Res.* **16**(14), 441–468 (2007)
8. Bordone, V.: Social norms and intergenerational relationships. In: De Santis, G. (ed.) *The Family, the Market or the State?*, pp. 159–178. Springer, Dordrecht (2012). doi:[10.1007/978-94-007-4339-7](https://doi.org/10.1007/978-94-007-4339-7)
9. Ludwig, A., Zimmerperz, A.: A parsimonious model of subjective life expectancy. WP 74, Econ. Res. S. Afr. (2008)
10. O’Brien, C., Fenn, P., Diacon, S.: How long do people expect to live? Results and implications. Centre for Risk and Insur. Stud. Res. Report 2005–1. Nottingham Un Business School (2005)
11. O’Connell, A.: How long do we expect to live? *Popul. Ageing* **4**, 185–201 (2011)
12. Berrington, A., Diamond, I.: Marital dissolution among the 1958 British birth cohort: the role of cohabitation. *Popul. Stud.* **53**(1), 19–38 (1999)
13. Impicciatore, R., Billari, F.: Secularization, union formation practices, and marital stability: evidence from Italy. *Eur. J. Popul.* **28**(2), online first (2012)
14. van Solinge, H., Henkens, K.: Living longer, working longer? The impact of subjective life expectancy on retirement intentions and behaviour. *Eur. J. Public Health* **20**(1), 47–51 (2010)
15. Margolis, R., Myrskylä, M.: A global perspective on happiness and fertility. *Popul. Dev. Rev.* **37**(1), 29–56 (2011)
16. de La Rochebrochard, É.: Assistance médicale à la procréation (AMP). *Dictionnaire de démographie* (dir. par F. Meslé, L. Toulemont et J. Véron) (2011)
17. De Rose, A., Salvini, S.: *Rapporto sulla popolazione*. Il Mulino, Bologna (2011)
18. Liefbroer, A.C.: Changes in family size intentions across young adulthood: a life-course perspective. *Eur. J. Popul.* **25**(4), 363–386 (2009)
19. Iacovou, M., Tavares, L.P.: Yearning, learning and conceding: reasons men and women change their childbearing intentions. *Popul. Dev. Rev.* **37**(1), 89–123 (2011)
20. Testa, M.R.: Family sizes in Europe: evidence from the 2011 Eurobarometer survey. *Eur. Demogr. Res. Paper 2* (2012)
21. UN: *Child Adoptions: Trends and Policies*. New York (2009)
22. Kohler, H.P., Behrman, J.R., Skytthe, A.: Partner + children = happiness? The effects of partnerships and fertility on well-being. *Popul. Dev. Rev.* **31**(3), 407–445 (2005)
23. Rivellini, G., Rosina, A., Sironi, E.: Marital disruption and subjective well-being: evidence from an Italian panel survey. Paper of the 46th SIS Scientific Meeting, Rome (2012)
24. Verbakel, E.: Subjective well-being by partnership status and its dependence on the normative climate. *Eur. J. Popul.* **28**(2), online first (2012)
25. Zimmermann, A.C., Easterlin, R.A.: Happily ever after? Cohabitation, marriage, divorce and happiness in Germany. *Popul. Dev. Rev.* **32**(3), 511–528 (2006)
26. Aassve, A., Betti, G., Mazzucco, S., Mencarini, L.: Marital disruption and economic well-being: a comparative analysis. *J. R. Stat. Soc. Series A* **170**(3), 781–799 (2007)
27. Piketty, T.: The impact of divorce on school performance : evidence from France, 1968–2002, CEPR discussion paper series no 4146 (2003)
28. Leridon, H.: La famille va-t-elle disparaître ? in *Dictionnaire de démographie* (dir. par F. Meslé, L. Toulemont et J. Véron) 164–166 (2011)
29. Istat: *La povertà in Italia*. Anno 2010 (2011)
30. OECD: *What are equivalence scales?* Paris (2008)
31. Istat: *La misurazione delle tipologie familiari nelle indagini di popolazione*. *Metodi e norme* 46 (2010)
32. Schizzerotto, A., Trivellato, U., Sartor, N.(a.c.d.): *Generazioni disuguali. Le condizioni di vita dei giovani di ieri e di oggi: un confronto*. Il Mulino, Bologna (2011)

STAR Modeling of Pulmonary Tuberculosis Delay-Time in Diagnosis

Bruno de Sousa, Dulce Gomes, Patrícia A. Filipe, Cristiana Areias, Teodoro Briz, Carlos Pires, and Carla Nunes

Abstract

Understanding what characterizes patients who suffer great delays in diagnosis of pulmonary tuberculosis is of great importance when establishing screening strategies to better control TB. Greater delays in diagnosis imply a higher chance for susceptible individuals to become infected by a *bacilliferous* patient. A structured additive regression model is attempted in this study in order to potentially contribute to a better characterization of *bacilliferous* prevalence in Portugal. The main findings suggest the existence of significant regional differences in Portugal, with the fact of being female and/or alcohol dependent contributing to an increased delay-time in diagnosis, while being dependent on intravenous drugs and/or being diagnosed with HIV are factors that increase the chance of an earlier diagnosis of pulmonary TB. A decrease in 2010 to 77 %

B. de Sousa (✉)

CINEICC, Faculdade de Psicologia e de Ciências da Educação – Universidade de Coimbra, Coimbra, Portugal
e-mail: bruno.desousa@fpce.uc.pt

D. Gomes • P.A. Filipe

CIMA/UE, Escola de Ciência e Tecnologia – Universidade de Évora, Évora, Portugal
e-mail: dmog@uevora.pt; pasf@uevora.pt

T. Briz • C. Areias

Escola Nacional de Saúde Pública – Universidade Nova de Lisboa, Lisboa, Portugal
e-mail: tshb@ensp.unl.pt; c.areias@ensp.unl.pt

C. Pires

Instituto de Higiene e de Medicina Tropical – Universidade Nova de Lisboa, Lisboa, Portugal
e-mail: carlosandrepirez@gmail.com

C. Nunes

CISP, Escola Nacional de Saúde Pública – Universidade Nova de Lisboa, Lisboa, Portugal
e-mail: CNunes@ensp.unl.pt

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,
Studies in Theoretical and Applied Statistics,
DOI 10.1007/978-3-319-27274-0_19

215

on treatment success in Portugal underlines the importance of conducting more research aimed at better TB control strategies.

1 Introduction

Although many studies have strongly indicated that tuberculosis can be controlled in almost any socio-economical reality [9, 14, 16], it remains a struggle to successfully control TB when faced with the presence of an epidemic HIV infection [8]. The Tuberculosis Programme from the European Center for Disease Prevention and Control (ECDC) recognizes that improvements have been made in tuberculosis (TB) prevention, but still considers it a threat to human health both world-wide and in Europe. Tuberculosis is currently classified as a re-emerging disease of European importance, with Portugal in 2008 still reporting a notification rates higher than 20 per 100,000 (21 per 100,000 in 2011 [3]), together with Romania, Lithuania, Latvia, Bulgaria, Estonia, and Poland [4]. In [3], Portugal, together with other European countries, has a sex ratio of men to women of 2:1, yet with a tendency to become more subtle in the future, as shown in other EU/EAA countries [5].

The HIV epidemic undermines the control of tuberculosis. Although the quality and completeness of country data on TB/HIV co-infection vary greatly, out of the eight countries that reported complete data in 2008, with co-infection rates between 0 and 14.6%, Portugal has one of the highest proportions of co-infections cases, together with Estonia and Malta. Nevertheless, Portugal, Iceland, and Slovakia achieved the target of a treatment success rate of 85% or higher set by the Stop TB Partnership. Among the 22 studied countries, the successful outcome among previously untreated culture-positive pulmonary TB cases in 2007 was 79.5% [4]. Unfortunately, latest data shows a decrease in 2010 to 77% on treatment success rate in Portugal, partly due to an increase in treatment time to 1 year [3].

Hornick in 2008 [10] reports that data from the Centers for Disease Control and Prevention (CDC) show that approximately 21–23% of individuals coming into close contact with patients suffering from pulmonary infectious tuberculosis themselves become infected. A multitude of factors could contribute to the increased risk of an individual contracting TB, as well as of disseminating it if already ill, pulmonary and contagious. With this in mind, the current study explores the effect of some of these factors on the delay-time in diagnosis of pulmonary TB in Portugal. The factors considered were the addictive consumption of alcohol, IV drugs (intravenous drugs), and other drugs, as well as sex, age, number of previous treatments, and HIV status of an individual. A structured additive regression (STAR) model was fit to the data in order to explore possible spatial correlations that can arise from the individual's municipality of residence together with the risk factors and other environmental variables, such as being an inmate, homeless or living in a risk area. Living in a risk area was a variable determined from a previously study by Nunes et al. [15] where geographical areas were classified as high/not high risk areas for contracting pulmonary TB.

STAR models [7] is the class of complex regression models chosen in this study since it allows to take into account a multitude of covariates while exploring possible spatial and temporal correlations. In particular, a structured hazard regression model is applied relaxing the strong condition of proportional hazards in the Cox model [2].

The material and methods are presented in Sect. 2, followed by the main results (Sect. 3), and a final discussion in Sect. 4.

2 Material and Methods

The database used was provided by two official sources, namely, the National Program for Tuberculosis Control (Pulmonary Tuberculosis notified cases between 2000 and 2009) and Statistics Portugal—INE (population data). The information provided include lifestyle characteristics of the individuals (alcohol, IV drugs or other drug dependence), characteristics inherent to the individual (sex, age, number of treatments, new case, HIV), and environmental variables (municipality of residence, being an inmate, homeless, living in a risk area). The delay-time variable will be the focus of our analysis in Sect. 3, representing the time between the first symptoms and the diagnosis of pulmonary tuberculosis. Table 1 contains a full description of all variables that will be pursued in the regression analysis performed in Sect. 4.

The results of our analysis will be based on $n = 13,615$ complete notified cases, i.e., the cases for which we have all the information regarding the variables defined in Table 1. This represents 58.7% of the original database with a delay-time less or equal to 365 days. By working only with complete data, it is worth noting that we observe similar percentages of *bacilliferous* and HIV patients (75.3%

Table 1 Description of variables of the Pulmonary Tuberculosis notified cases database

Variable	Description
DelayTime	Time between the first symptoms and the diagnosis of tuberculosis
Municipality	Municipality where the individual lives (278 municipalities, excluding the autonomous regions of Azores and Madeira)
Sex	Gender of the individual with categories “male” (= 0) and “female” (= 1)
Age	Age of the individual in years
Ntreatments	Number of treatments before present diagnosis
Alcohol	Whether an individual is alcohol dependent (1 = Yes and 0 = No)
IVDrugs	Whether an individual is dependent on IV drugs (1 = Yes and 0 = No)
OtherDrugs	Whether an individual is dependent on other drugs (1 = Yes and 0 = No)
Inmate	Whether an individual is an inmate (1 = Yes and 0 = No)
Homeless	Whether an individual is homeless (1 = Yes and 0 = No)
HIV	Whether an individual has HIV (1 = Yes and 0 = No)
RiskArea	Whether an individual lives in a risk area (1 = Yes and 0 = No)
NewCase	Whether an individual is a new diagnosed case (1 = Yes and 0 = No)

and 14.8 %, respectively) when compared to official data reports (70–75 % and 13–15 %, respectively [3]). With this in mind, a brief descriptive comparison analysis was performed between the data analyzed in this study and the notified cases that were omitted due to missing values. Among the three variables for which we had information for all the notified cases, namely, age, delay-time in diagnosis, and sex, the differences between the two groups were very slight, with both groups following the same patterns of behavior.

Survival time and censoring can be modeled through the Cox proportional hazards model proposed by David Cox in 1972 [2]. In this model the hazard rate, $\lambda(t|u) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} P(t \leq T < t + \Delta t | T \geq t, u)$, can be interpreted as the instantaneous rate of an event in the interval $[t, t + \Delta t]$, given survival up to time t .

The main goal of survival regression is to describe the influences of covariates u through a regression model for the hazard rate. In the Cox proportional hazards model, the hazard rate is assumed to have a multiplicative structured of the form

$$\lambda(t|u) = \lambda_0(t)\exp(u'\gamma), \quad (1)$$

where $\lambda_0(t)$ is an unspecific baseline hazard rate and $u'\gamma$ is a linear predictor formed of (time-constant) covariates u and regression coefficients γ .

Because the ratio between the hazard rates for two individuals with covariates vectors u_1 and u_2 is independent of t , the Cox model (1) is called a proportional hazards model. This assumption is not always present in real life data and needs to be checked for the model to be applied. To account for nonproportional hazards, nonstandard covariate effects, and spatial dimension, the classical Cox model is extended to a nonparametric structured hazard rate model [11] defined as $\lambda_i(t) = \exp(\eta_i(t))$, $i = 1, \dots, n$, with the structured additive predictor defined as:

$$\eta_i(t) = u_i(t)'\gamma + g_0(t) + \sum_{k=1}^K g_k(t)w_{ik}(t) + \sum_{j=1}^J f_j(v_{ij}(t)), \quad (2)$$

where $g_0(t) = \log(\lambda_0(t))$ is the log-baseline hazard, $g_k(t)$ represents time-varying effects of covariates $w_{ik}(t)$, $f_j(v_{ij}(t))$ are non-linear effects of different types of generic covariates and $u_i(t)'\gamma$ corresponds to effects of parametric covariates.

In this study, a geoadditive predictor is chosen for the hazard rate $\lambda(t) = \exp(\eta(t))$, with $\eta(t)$ defined as:

$$\eta(t) = g_0(t) + f_1(\text{Age}) + f_2(\text{Ntreat}) + f_{\text{str}}(\text{Munic}) + f_{\text{unstr}}(\text{Munic}) + u(t)'\gamma, \quad (3)$$

where $g_0(t)$ denotes the log-baseline hazard rate, f_1, f_2 are functions of the covariates age and number of treatments. Functions f_{str} and f_{unstr} model the global and local spatial effects based on the municipality where an individual lives, representing the effects that obey a strong spatial structured and the ones that are present locally,

respectively. The fixed effects of the numerous categorical covariates (Table 1) are represented by $u(t)$.

Inference is based on a mixed model representation of the structured additive predictor (2) and yields either penalized likelihood estimates (from a frequentist perspective) or empirical Bayes/posterior mode estimates (from a Bayesian perspective) [12]. The structured additive predictor (2) can be partitioned into an unpenalized part (fixed effects) and a penalized part (random effects), where a flat prior and an i.i.d. Gaussian prior probability distributions are considered, respectively. Both effects are estimated using REML-restricted maximum likelihood [7, 11]. In particular, the non-linear effects, f_j , are estimated based on Bayesian P-splines [13]; and the spatial effects follow a Markov random field approach where two municipalities are considered neighbors if they share a common boundary and the effect of the municipality is conditionally Gaussian [6, 12].

The estimation of the models was performed in the open source software packages BayesX [1] and R [17].

3 Results

The estimates for the log-baseline g_0 and the nonparametric effects f_j from model (3) are shown in the next figure. The log-baseline, Fig. 1a, shows a steep increase until approximately 44 days, followed by an alternate period between positive and negative effects (approximately constant) until around 340 days. At the end of the observation period, there is a strong increase in g_0 . However, only 44 individuals had a delay-time in diagnosis more than 340 days and, therefore, this increase should not be over-interpreted.

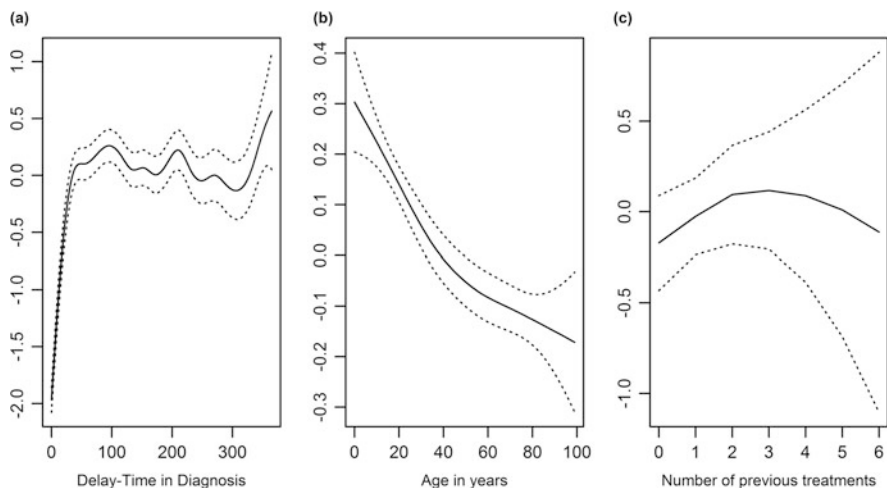


Fig. 1 Municipality-level analysis: posterior mode estimates of the effects of the log-baseline (a), age (b), and the number of previous treatments (c), together with pointwise 95% credible intervals (dashed lines)

From Fig. 1 there is a non-linear effect for both effects f_1 and f_2 in Eq. (3), regarding age and number of previous treatments before the current diagnosis, with a clear stronger effect in the case of the age effect. The chance of the event decreases with age, meaning that younger people have a higher chance of being diagnosed earlier, with a clear decrease as people grow older (Fig. 1b). Worth noting that after 45 years of age this decrease slows down when compared with younger people.

The effect of the number of previous treatments before the current diagnosis is almost constant with a slight increase for those who had at least two treatments before. However, since the credible intervals include zero, the influence of the number of previous treatments can be neglected (Fig. 1c). Also of note is the large bandwidth for values higher than 2. This is due to the small number of cases of people with more than 2 previous treatments (44 cases).

Looking at the estimated global spatial effects in the left panel of Fig. 2a, we find that municipalities with higher chances of a delayed diagnosis seem to occur at the

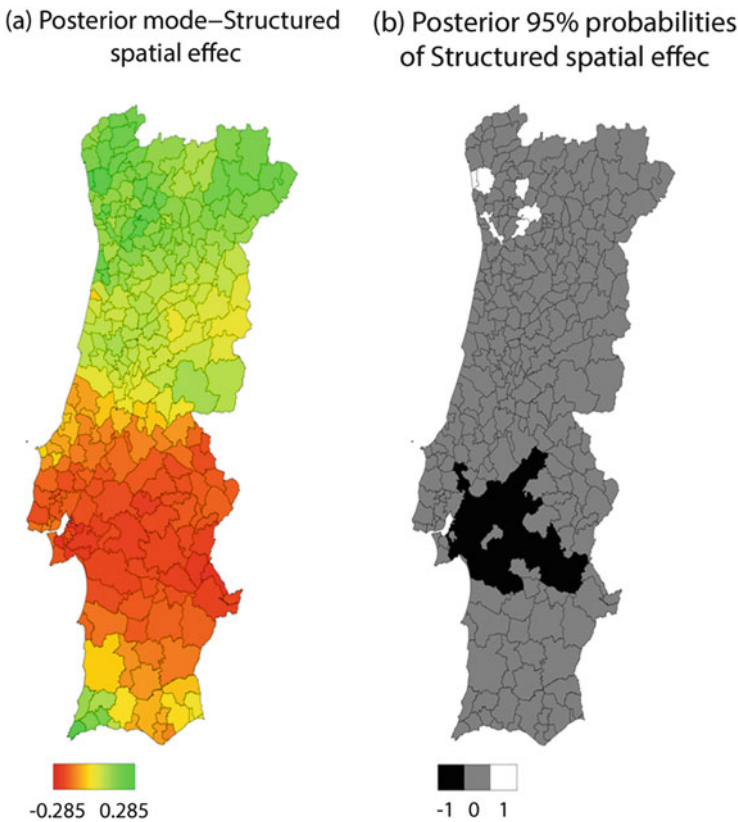


Fig. 2 Municipality-level analysis: Estimated global (a) spatial effects and pointwise 95 % significance map (b). *Black* denotes municipalities with strictly negative credible intervals, whereas *white* denotes municipalities with strictly positive credible intervals

center and south of Portugal (red regions). The statistical significance of this areas can be confirmed by the map of Fig. 2b, where black denotes municipalities with strictly negative credible intervals and white denotes municipalities with strictly positive credible intervals, i.e., representing a gray scale of municipalities that clearly contribute to an increased chance of a delayed diagnosis (in black) as compared to the ones which contribute to decreasing that chance (in white). Worth noting that the Lisbon area and upper Alentejo (black areas, Fig. 2b) are the regions contributing to an increase in the delay, while a decreasing effect in the delay is present in a few regions in the Oporto area (white areas, Fig. 2b).

In terms of the estimated local spatial effects in Fig. 3a, we observe quite a homogeneous map with most of the regions grey and some areas pink and green, indicating no local effects in general, with some municipalities showing a tendency for delayed diagnosis. This structure is confirmed by the significance map in Fig. 3b,

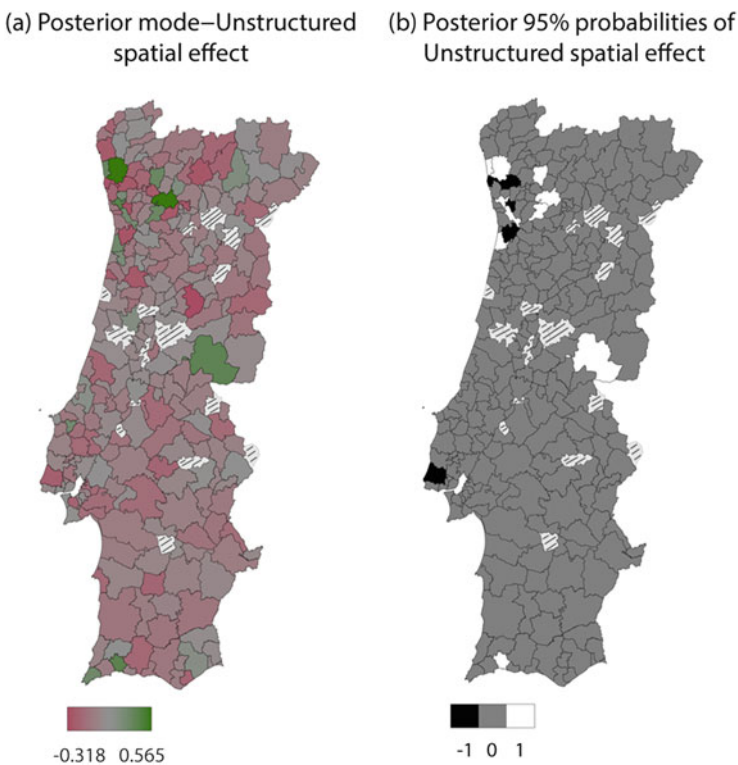


Fig. 3 Municipality-level analysis: estimated local (a) spatial effects and pointwise 95 % significance map (b). *Black* denotes municipalities with strictly negative credible intervals, whereas *white* denotes municipalities with strictly positive credible intervals. Regions without observations are represented with *diagonal stripes*

Table 2 Municipality-level analysis: estimates (Post. Modes), standard deviations, p -values, and 95 % credible intervals of fixed effects

Variable	Estimates	Std. dev.	p -value	95 % credible interval	
const	-4.218	0.138	< 0.001	-4.488	-3.947
Sex	-0.103	0.020	< 0.001	-0.142	-0.063
Alcohol	-0.068	0.024	0.006	-0.116	-0.020
IVDrugs	0.109	0.039	0.006	0.032	0.186
OtherDrugs	-0.036	0.035	0.310	-0.105	0.033
Inmate	0.107	0.072	0.137	-0.034	0.248
Homeless	0.064	0.068	0.341	-0.068	0.197
HIV	0.170	0.029	< 0.001	0.112	0.227
RiskArea	0.058	0.058	0.317	-0.056	0.173
NewCase	0.096	0.182	0.596	-0.261	0.453

where black denotes districts with strictly negative credible intervals (higher chance of a delayed diagnosis) and white denotes districts with strictly positive credible intervals. The municipalities with an increased local chance for a delayed diagnosis are in the Oporto and Lisbon areas.

Table 2 contains the estimates, standard deviations, p -values, and 95 % credible intervals of the fixed effects of model in Eq. (3). Being an individual in a risk area, a new case, an inmate, a homeless person, or one dependent on other drugs does not seem to be a statistically significant factor in the delay-time in diagnosis. Among the statistically significant factors, being female and/or alcohol dependent seems to reduce the chance of the event (being diagnosed), i.e., increasing the delay-time in diagnosis. On the other hand, being dependent on IV drugs and/or being diagnosed with HIV increase the chances of an earlier diagnosis of TB.

In addition, we conducted a district-level analysis where the 278 municipalities were classified into 18 districts. The covariates and fixed effects were very similar to the results presented above. Of note was the very clear pattern in terms of spatial effects presented in Fig. 4. Districts in the center and south of Portugal seem to have a higher chance of a delayed diagnosis. In terms of the estimated local spatial effects in the right panel of Fig. 4, it is very clear that Lisbon and Oporto are the districts that emerged as those that contribute to an increase or decrease in the delay-time in diagnosis, respectively. Neither global nor local spatial effects lead to strictly negative or positive credible intervals.

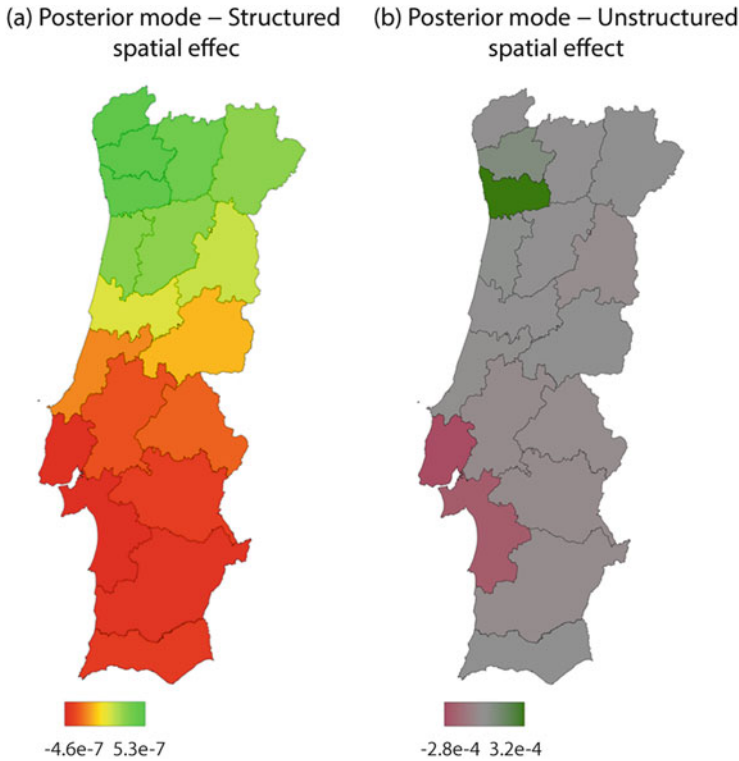


Fig. 4 District-level analysis: estimated global (a) and local (b) spatial effects

4 Discussion and Conclusions

The delay-time is the time between the appearance of the first symptoms and the diagnosis of a pulmonary TB case by the health care system. This time period depends on the actions of both the patient and health care. In Portugal, 95 % of the TB cases are diagnosed because symptomatic (with cough) patients search health care services [3]. This means that when a patient is diagnosed, he/she may well have already infected someone, which can lead to an endless endemic state.

Understanding what characterizes those patients who suffer great delays in diagnosis may contribute, for example, to establishing a better detection, and therefore, a decrease of the endemic level in Portugal. Our study suggests that younger people have an higher chance of being diagnosed earlier, with a clear decrease as people grow older. As reported in [3], age patterns have been changing among TB patients in Portugal as well as in other developed countries. Over time,

a more pronounced decrease of incidence is observed in younger individuals (0–44 years) as opposed to older groups (especially 45–74 years) [3]. Future analysis will focus on these age groups in order to better understand the possible reasons for this behavior.

Municipalities with longer delay-time in diagnosis seem to occur in the center and south regions of Portugal, with the Lisbon area and upper Alentejo being the regions which significantly contribute to an increase in the delay. Further research is needed in order to identify possible reasons that go beyond the factors addressed in this study that might affect these municipalities and explain the results obtained. Three important factors that can clearly contribute to a delay in the diagnosis are: (1) the population's lack of knowledge, (2) inefficient health care services, and (3) low incidence of the disease, implying a lesser chance of it being a primary diagnosis.

Among the factors that were considered in this study, being female and/or being alcohol dependent indicates a tendency for an increasing delay-time in diagnosis, while being dependent on IV drugs and/or being diagnosed with HIV increases the chance of an earlier diagnosis of pulmonary TB. The knowledge of HIV status has been increasing in Portugal, from 59 % of missing cases in 2000 to 41 % in 2009. In future research the evolution of the results according to year of notification will be considered in order to be able to identify possible biases due to missing data.

With this initial study, we attempted to understand some of the main risk factors that might be responsible for greater delays in diagnosis, a pivotal piece of the puzzle in order to successfully control TB. Further studies are needed in order to clearly understand the problem within a more global perspective and address some of the research questions stated in this discussion.

Acknowledgements This work was financed by FCT/MCTES through the research project PTDC/SAU-SAP/116950/2010.

References

1. Belitz, C., Brezger, A., Kneib, T., Lang, S., Umlauf, N.: BayesX - software for Bayesian inference in structured additive regression models. Version 2.0. <http://www.bayesx.org/> (2009)
2. Cox, D.R.: Regression models and life tables (with discussion). *J. R. Stat. Soc. B* **34**, 187–220 (1972)
3. DGS: Relatório StopTb2012-Ponto da Situação Epidemiológica e de Desempenho (2012). Available in http://www.portaldasaude.pt/NR/rdonlyres/8E0DFF04-F030-43B4-80EB-A71AD96F3718/0/relatorio_tuberculose_2012.pdf
4. ECDC: Progressing towards TB elimination – a follow-up to the Framework Action Plan to Fight Tuberculosis in the European Union. ECDC Special Report (2010). Available in http://www.ecdc.europa.eu/en/publications/Publications/101111_SPR_Progressing_towards_TB_elimination.pdf
5. European Centre for Disease Prevention and Control/WHO Regional Office for Europe. Tuberculosis surveillance and monitoring in Europe 2012. Stockholm: European Centre for Disease Prevention and Control (2012). Available in <http://ecdc.europa.eu/en/publications/Publications/1203-Annual-TB-Report.pdf>
6. Fahrmeir, L., Lang, S.: Bayesian inference for generalized additive mixed models based on Markov random fields priors. *J. R. Stat. Soc. C* **50**, 201–220 (2001)

7. Fahrmeir, L., Kneib, T., Lang, S.: Penalized structured additive regression for space-time data: a Bayesian perspective. *Stat. Sin.* **14**, 731–761 (2004)
8. Frieden, T.R.: Can tuberculosis be controlled? *Int. J. Epidemiol.* **31**, 894–899 (2002)
9. Frieden, T.R., Fujiwara, P.I., Washko, R.M., Hamburg, M.A.: Tuberculosis in New York City—turning the tide. *New Engl. J. Med.* **333**, 229–233 (1995)
10. Hornick, D.B.: Tuberculosis. In: Wallace, R.B. (ed.) *Maxcy-Rosenau-Last Public Health and Preventive Medicine*, pp. 248–257, 15th edn. McGraw-Hill Medical, New York (2008)
11. Kneib, T.: *Mixed Model Based Inference in Structured Additive Regression*. Fakultät für Mathematik, Informatik und Statistik der Ludwig-Maximilians-Universität München, Germany (2006)
12. Kneib, T., Fahrmeir, L.: A mixed model approach to structured hazard regression. *Scand. J. Stat.* **34**, 207–228 (2007)
13. Lang, S., Brezger, A.: Bayesian P-splines. *J. Comput. Graph. Stat.* **13**, 183–212 (2004)
14. Marrero, A., Caminero, J.A., Rodriguez, R., Billo, N.E.: Towards elimination of tuberculosis in a low-income country: the experience of Cuba, 1962–97. *Thorax* **55**, 39–45 (2000)
15. Nunes, C., Briz, T., Gomes, D., Filipe, P.A.: Pulmonary tuberculosis and HIV/AIDS: joint space-time clustering under an epidemiological perspective. In: Cafarelli, B. (eds.) *Proceedings of the Spatial Data Methods for Environmental and Ecological Processes - 2nd Edition*, 1–4, Foggia e Gargano (2011)
16. Suarez, P.G., Watt, C.J., Alarcon, E., Portocarrero, J., Zavala, D., Canales, R., Luelmo, F., Espinal, M.A., Dye, C.: The dynamics of tuberculosis in response to 10 years of intensive control effort in Peru International. *J. Infect. Dis.* **184**, 473–478 (2001)
17. R Core Team: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna (2012). <http://www.R-project.org>

Non-aggregative Assessment of Subjective Well-Being

Marco Fattore, Filomena Maggino, and Alberto Arcagni

Abstract

In this paper, we introduce a new methodology for socio-economic evaluation with ordinal data, which allows to compute synthetic indicators without variable aggregation, overcoming some of the major problems when classical evaluation procedures are employed in an ordinal setting. In the paper, we describe the methodology step by step, discussing its conceptual and analytical structure. For exemplification purposes, we apply the methodology to real data pertaining to subjective well-being in Italy, for year 2010.

1 Introduction

The use of ordinal data is spreading in socio-economic analysis, as issues like evaluating multidimensional poverty, well-being and quality-of-life are gaining importance in applied research and policy-making. Many social surveys ask respondents for self-assessments or subjective judgments, often expressed through binary or ordinal scales. Nowadays, many datasets comprising (also) ordinal variables are available to scholars; main examples at European and Italian level are the EU-SILC survey and the multi-topic survey entitled “Multipurpose Survey about Families Aspects of Daily Life”, held by Istat (Italian National Statistical Bureau). Despite the abundance of ordinal data, statistical methodologies capable to effectively

M. Fattore (✉) • A. Arcagni

Department of Statistics and Quantitative Methods, University of Milano - Bicocca, Milano, Italy
e-mail: marco.fattore@unimib.it

F. Maggino

Department of Statistics, Computer Science, Applications “G. Parenti”, University of Florence, Firenze, Italy

exploit them in studies pertaining to socio-economic evaluation are missing yet. Often, ordinal scores are treated as or transformed into cardinal figures and standard multivariate procedures are applied. Alternatively, ordinal data are simplified into binary variables and counting procedures are employed [1]. In both cases, the informative potential of ordinal data is not adequately exploited. As a matter of fact, the development of evaluation procedures in multidimensional ordinal settings is still an open, and largely unexplored, research field. Recently, a new methodology has been proposed by the Authors and other colleagues with the aim of overcoming counting and composite indicators approaches [4–6]. Its novelty lies in the use of partial order theory as a tool to compute synthetic indicators without aggregating ordinal variables. In the following, we give a step-by-step description of the methodology and, for exemplification purposes, apply it to data pertaining to subjective well-being in Italy. The paper is organized as follows; Sect. 2 introduces the data and motivates the interest for subjective well-being; Sect. 3 outlines the methodology; Sect. 4 presents the results of the analysis; Sect. 5 concludes.

2 Subjective Well-Being Data

The measurement of well-being is one of the most vivid topics in socio-economic statistics, particularly after the Stiglitz Commission stated a growing role of well-being measures, besides GDP, to assess the wealth of countries. In Italy, an ambitious project devoted to well-being assessment is being led by CNEL (National Committee for Economy and Work) and Istat. Twelve well-being dimensions have been identified; among them, our focus is on subjective well-being. Data used in the paper come from “Multipurpose Survey about Families Aspects of Daily Life” for year 2010.¹ The sample is composed of 48,336 statistical units. For sake of simplicity, records with missing values have been deleted reducing the sample to 40,949 units (see Sect. 4 for a remark on the missing data problem). We have selected four variables, pertaining to the satisfaction degree relative to:

1. personal economic status (variable v_1);
2. personal health status (variable v_2);
3. relationships with relatives (variable v_3);
4. relationships with friends (variable v_4).

All of the variables are recorded on a four-degree scale² (1 = “not at all”; 2 = “not much”; 3 = “enough”; 4 = “very”). In addition to well-being scores,

¹Data are available within a protocol agreement signed by Istat and the University of Florence.

²In the original dataset, variables are scored as: 4 = “not at all”; 3 = “not much”; 2 = “enough”; 1 = “very”. Codes have been reversed in such a way that increasing scores correspond to increasing satisfaction.

information about gender and the region of residence of each statistical unit in the sample are available.

3 Evaluating Subjective Well-Being

As usual in studies pertaining to well-being, the primary aim is two-fold: (1) identifying people who are not satisfied of their own well-being status and (2) measuring their dissatisfaction degree. The existence of incomparabilities among well-being self-assessments makes these goals more subtle than in the unidimensional case, where individual achievements can be linearly ordered. Multidimensional self-assessments can be ordered only partially and this introduces the role of partial order theory in the evaluation procedure.

The Partial Order of Well-Being Self-assessments By self-assessments, any statistical unit in the population is assigned a four-component vector \mathbf{p} , in the following called a *profile*, comprising his/her scores on variables v_1, v_2, v_3 and v_4 . In total, there are $4^4 = 256$ different profiles $\mathbf{p}_1, \dots, \mathbf{p}_{256}$, together with the (absolute) frequencies n_1, \dots, n_{256} of statistical units sharing them. Profiles can be partially ordered according to the following natural definition:

Definition 3.1 Profile \mathbf{p}_h is more satisfied than, or equally satisfied as, profile \mathbf{p}_k (written $\mathbf{p}_k \preceq \mathbf{p}_h$) if and only if $p_{ki} \leq p_{hi}$ for each $i = 1, \dots, 4$, where p_{hi} and p_{ki} are the i -th components of \mathbf{p}_h and \mathbf{p}_k , respectively.

The set P of profiles endowed with the partial order \preceq gives rise to the *profile poset* (P, \preceq) , which, for notational convenience, will be similarly indicated as P . It has a top element (4444), denoted by \top , and a bottom element (1111), denoted by \perp , which represent the best and the worse element, respectively.

Setting the Threshold Given P , the open problem is how to extract information pertaining to well-being, out of it. The identification of unsatisfied profiles (just like the identification of poor individuals in customary poverty studies) is a normative act, which cannot be performed only through data analysis. As a purely mathematical structure, P conveys no absolute socio-economic information and cannot suffice to identify unsatisfied profiles. Identification is therefore performed introducing exogenously a *threshold* (here denoted by τ), which in principle is up to experts and policy-makers to select. In the literature about social evaluation, multidimensional thresholds are usually identified based on the selection of cut-offs for each evaluation dimension. In a partial order framework, where emphasis is put on profiles rather than variables, it is more natural to identify the threshold directly in terms of profiles on “the edge of dissatisfaction”. This way, one can take into account interactions among achievements on different well-being factors, which are crucial in a multidimensional setting. We must also notice that due to multidimensionality, more than one profile may be on the dissatisfaction edge and thus the threshold τ

may be (and usually is) composed of several elements. As proved in [4], under very general conditions, τ may be always chosen as an *antichain* of P , that is, as a set of mutually incomparable elements of the profile poset. It is clear that, in real studies, the choice of the threshold is a critical task, affecting all the subsequent results. Therefore preliminary data insights, experts' judgments and any other source of information should be involved in selecting it.

Identification of Unsatisfied Profiles Given the threshold, the next step is to define an *identification function*, denoted by $\text{idn}(\cdot)$, that quantifies in $[0, 1]$ to what extent a profile of P may be classified as unsatisfied. Notice that $\text{idn}(\cdot)$ does not measure the intensity of dissatisfaction (which will be later assessed in a different way), but the degree of membership to the set of unsatisfied profiles. The methodology is thus fuzzy in spirit, to reflect the classification ambiguities due to multidimensionality and partial ordering. In view of its formal definition, it is natural to impose the following four conditions on $\text{idn}(\cdot)$:

1. If, in satisfaction terms, profile p is better than profile q , then its degree of membership to the set of unsatisfied profiles must be lower than the degree of q , in formulas:

$$q \preceq p \Rightarrow \text{idn}(p) \leq \text{idn}(q).$$

2. Profiles belonging to the threshold are, by definition, unsatisfied profiles; therefore the identification function must assume value 1 on them:

$$p \in \tau \Rightarrow \text{idn}(p) = 1.$$

From conditions (1) and (2), it follows that a profile is unambiguously classified as unsatisfied if it belongs to the threshold or if it is worse than an element of the threshold:

$$\text{idn}(p) = 1 \Leftrightarrow p \preceq q, \quad q \in \tau.$$

In poset theoretical terms, the subset of profiles satisfying the above condition is called the *downset* of τ and is denoted by $\tau \downarrow$.

3. A profile p is unambiguously identified as “not unsatisfied” if and only if it is better than *any* profile belonging to the threshold:

$$\text{idn}(p) = 0 \Leftrightarrow q \preceq p, \quad \forall q \in \tau.$$

This condition aims to exclude that a profile which is incomparable with even a single element of the threshold may be scored 0 by the identification function.

4. If P is a linear order, then $\text{idn}(\cdot)$ must assume only values 0 or 1. In other words, if no incomparability exists, the identification function must classify profiles as either unsatisfied or not.

To determine the functional form of the identification function, we start by defining $\text{idn}(\cdot)$ on the linear extensions of the profile poset. A *linear extension* ℓ of P is a linear order defined on the set of profiles and obtained turning incomparabilities of P into comparabilities. In a linear extension, all the elements are comparable, particularly the elements of the threshold τ selected in P . Therefore, in any linear extension ℓ , we may find an element τ_ℓ of the threshold that is ranked above any other element of τ . According to condition (4), it is natural to define the identification function³ on ℓ putting $\text{idn}_\ell(\mathbf{p}) = 1$ if $\mathbf{p} \preceq_\ell \tau_\ell$ ⁴ and $\text{idn}_\ell(\mathbf{p}) = 0$ otherwise. In other words, identification in linear extensions reduces to the unidimensional problem of classifying profiles as above the threshold or not. We now extend the definition of the identification function from the set of linear extensions to the profile poset. The starting point is a simple but fundamental results of partial order theory that we state without proof [7]:

Theorem 3.1 *Any finite poset P is the intersection of its linear extensions:*

$$P = \bigcap_{\ell \in \Omega(P)} \ell$$

where $\Omega(P)$ is the set of linear extensions of P .

The close connection between P and $\Omega(P)$ suggests to express $\text{idn}(\cdot)$ as a function of the idn_ℓ s:

$$\text{idn}(\cdot) = F(\{\text{idn}_\ell(\cdot), \ell \in \Omega(P)\}).$$

To specify the functional form of $F(\cdot, \dots, \cdot)$, we require it (1) to be symmetric (the way linear extensions are listed by is unimportant) and to satisfy the properties of (2) associativity, (3) monotonicity, (4) homogeneity and (5) invariance under translations. Symmetry and associativity are justified since the intersection operator is symmetric and associative; monotonicity assures that $\text{idn}(\mathbf{p})$ increases as the number of linear extensions where $\text{idn}_\ell(\mathbf{p}) = 1$ increases; homogeneity and invariance under translations assure that $\text{idn}(\cdot)$ changes consistently if the $\text{idn}_\ell(\cdot)$ s are rescaled or shifted. By the theorem of Kolmogorov–Nagumo–de Finetti, it then follows that $F(\cdot, \dots, \cdot)$ has the form of an arithmetic mean,⁵ so that:

$$\text{idn}(\mathbf{p}) = \frac{1}{|\Omega(P)|} \sum_{\ell \in \Omega(P)} \text{idn}_\ell(\mathbf{p}).$$

³We denote this identification function by idn_ℓ to remind that it depends upon the linear extension considered.

⁴We denote with \preceq_ℓ the order relation in ℓ .

⁵More precisely, of a weighted arithmetic mean, but in our case there is no reason to assign different weights to different linear extensions.

Since idn_ℓ is either 0 or 1, $\text{idn}(\mathbf{p})$ may be alternatively seen as the fraction of linear extensions where \mathbf{p} is classified as unsatisfied:

$$\text{idn}(\mathbf{p}) = \frac{|\{\ell \in \Omega(P) : \text{idn}_\ell(\mathbf{p}) = 1\}|}{|\Omega(P)|}.$$

In a sense, the evaluation procedure implements a counting approach, but on linear extensions of P and not directly on well-being variables. This way, it exploits the structure of the underlying partial order, to quantify the degree of membership of a profile to the set of unsatisfied profiles, with *no* variable aggregation. By construction, all of the elements in $\tau \downarrow$ are classified as unsatisfied in any linear extension of P and therefore are scored to 1 by $\text{idn}(\cdot)$, as required by condition (2) above. Similarly, profiles above any element of τ are scored to 0, consistently with condition (3). All of the other profiles in P are classified as unsatisfied in some linear extensions and as not unsatisfied in others and thus are scored in $]0, 1[$ by $\text{idn}(\cdot)$. Once each profile \mathbf{p} (and thus any statistical unit sharing it) has been scored by $\text{idn}(\cdot)$, synthetic well-being indicators may be obtained. In particular, we focus on the (fuzzy extension of the) *Head Count Ratio*, which is defined as the arithmetic mean of the identification function over the entire population and which represents the “relative amount” of dissatisfaction in it.

Measuring Dissatisfaction Intensity Two profiles may share the same identification degree, but still represent conditions of different dissatisfaction severity. Consider, for example, profile (4144), which belongs to the threshold, and profile (1111), which is the bottom of P , both scored 1 by the identification function. To obtain a more complete picture of subjective well-being, it is therefore of interest to separately assess the dissatisfaction intensity of a profile \mathbf{p} , which in the following will be called the *gap*⁶ of \mathbf{p} . To this goal, we:

1. Introduce a metric $d(\cdot, \cdot)$ on linear orders, to measure the distance between a profile and the threshold in each linear extension ℓ of P .
2. Given a linear extension ℓ , for any profile \mathbf{p} classified as unsatisfied in it, its distance $d(\mathbf{p}, \tau | \ell)$ to the threshold is computed. This distance is then scaled to $[0, 1]$, dividing it by the maximum distance to the threshold achievable in ℓ , that is, by $d(\perp, \tau | \ell)$. The rescaled distance is denoted by $\hat{d}(\mathbf{p}, \tau | \ell)$.
3. Similarly to the identification step, the gap $g(\mathbf{p})$ of profile \mathbf{p} is obtained averaging distances $\hat{d}(\mathbf{p}, \tau | \ell)$ over the set of linear extensions $\Omega(P)$.

Many different metrics may be defined on a linear extension. Here we simply define it as the absolute value of the difference between the rank of a profile and the rank of the highest ranked element of the threshold. Formally, let $r(\mathbf{p} | \ell)$ be the rank

⁶The terminology is taken by the practice of poverty measurement.

of profile \mathbf{p} in linear extension ℓ and let

$$r(\tau \mid \ell) = \max_{\mathbf{q} \in \tau} (r(\mathbf{q} \mid \ell))$$

be the rank of the highest ranked element of the threshold in ℓ . Then the distance between a profile and the threshold is simply $d(\mathbf{p}, \tau; \ell) = |r(\mathbf{p} \mid \ell) - r(\tau \mid \ell)|$ and, since $d(\perp, \tau; \ell) = r(\tau \mid \ell) - 1$, we also have

$$\hat{d}(\mathbf{p}, \tau; \ell) = \frac{|r(\mathbf{p} \mid \ell) - r(\tau \mid \ell)|}{r(\tau \mid \ell) - 1}.$$

Finally we put

$$g(\mathbf{p}) = \frac{1}{|\Omega(P)|} \sum_{\ell \in \Omega(P)} \hat{d}(\mathbf{p}, \tau \mid \ell).$$

Some comments on the gap function are in order. First, it is computed only for profiles \mathbf{p} such that $\text{idn}(\mathbf{p}) > 0$. Secondly, it is anti-monotonic, since clearly if $\mathbf{q} \neq \mathbf{p}$ and $\mathbf{q} \preceq \mathbf{p}$, then $d(\mathbf{p}, \tau \mid \ell) < d(\mathbf{q}, \tau \mid \ell)$ in each linear extension ℓ and thus $g(\mathbf{p}) < g(\mathbf{q})$. Thirdly, the gap function achieves its maximum value 1 on the bottom element of P (in our case, on profile (1111)). In general, it attains strictly positive values even on the elements of the threshold, achieving value 0 if and only if the threshold is composed of a single profile. This fact, due to the existence of incomparabilities among elements of the threshold, reveals how subtle multidimensional evaluation may be, compared to the unidimensional case. Once the gap function is computed, it may be averaged on the entire population, to obtain the overall *Gap* indicator which complements the Head Count Ratio previously introduced.

Computational Aspects The number of linear extensions of a poset like that involved in the present paper is too huge to list them and perform exact calculations of the identification and the gap functions. In practice, one extracts a sample of linear extensions and computes approximate results on it. The most effective algorithm for (quasi) uniform sampling of linear extensions is the Bublely–Dyer algorithm [2]. For the purposes of this paper, the algorithm has been implemented through a C routine, which is part of an R [8] package for poset calculations, under development by the Authors [3]. The computations required the extraction of 10^{10} linear extensions and took approximately 7.5 h on a 1.9 GB Intel Core 2 Duo CPU E8400 3.00 GHz \times 2, with Linux Ubuntu 12.04 64 bit.

4 Application to Subjective Well-Being Data

To show the evaluation methodology in action, we now apply it to the data presented in Sect. 2. First of all, a threshold has been selected, namely the antichain $\tau = \{(1144), (1211), (3111)\}$. More emphasis has been given to dissatisfaction relative to economics (first component of the profiles) and health (second component of the profiles), than to dissatisfaction pertaining to relationships with friends and relatives. However, it is the combination of scores that matters in identifying unsatisfied profiles. Consider, for example, the first element of the threshold: no matter how good relationships with friends and relatives are, if an individual reports heavy economic and health problems, the corresponding profile will be scored 1 by the identification function. Similarly, if the health status is slightly better, but relational problems arise, then the profile is again scored to 1 by the evaluation function (second element of the threshold). Analogously, for the economic dimension. The choice of the threshold requires in fact judgments on the “global meaning” of the profiles.⁷ Compensations among dimensions may exist, but this may depend upon the achievement levels in complex ways. Our choices could be argued indeed, but what is relevant here is to consider the flexibility of the approach, which allows to tune the threshold according to the aims and the contexts. Chosen the threshold, the identification function has been computed. The result is reported in Fig. 1. As it may be seen, its values range from 0 to 1 and some “levels” may be identified. Some profiles are almost unsatisfied, other are “just a little” unsatisfied and so on. This shows how the proposed procedure is successful in revealing the nuances of subjective well-being, overcoming rigid black or white classifications. A similar computation has been performed to get the gap function. Table 1 reports the results at regional and national level, also split by males and females. The Head Count Ratio ranges from about 7 % to almost 25 % and Gap ranges from about 8 % to about 16 %, revealing heavy interregional differences. Regions clearly separate in three main groups, below, around or above the national levels for both indicators. Broadly speaking, this distinction reflects the North–South axis, which is a typical feature of the Italian socio-economic setting, where southern regions are generally in worse socio-economic situations than the northern ones. However there is some remarkable shuffling among territorial areas and some regions from the South (Molise and Basilicata) turn out to score similarly to regions from the Centre and *vice versa*, as in the case of Umbria. The position of Trentino-Alto Adige is remarkable and confirms that this region is an outlier in the Italian context, due to its prerogatives and autonomy as a region under special statute and thanks to the efficiency of its administrative system. A closer look to Table 1 reveals also

⁷The choice of the threshold requires exogenous judgments and assumptions by social scientists and/or policy-makers. It must be noted, however, that the methodology allows for such exogenous information to be introduced in the analysis in a neat and consistent way. One could also add to the analysis judgments on the different relevance of well-being dimensions. Partial order theory, in fact, provides the tools to handle this information in a formal and effective way. We cannot give the details here, but some hints can be found in [6].

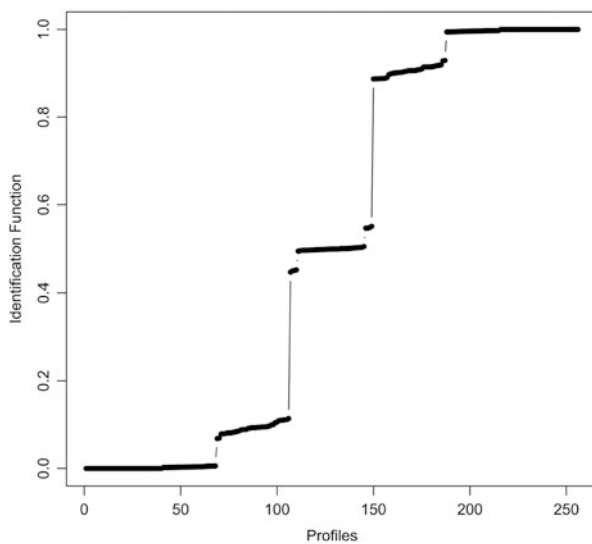


Fig. 1 Identification function (profiles ordered by increasing values of the identification function)

Table 1 Head Count Ratio and Gap at regional and national level

Regions	ID	Head count ratio			Gap		
		Total	Males	Females	Total	Males	Females
Piemonte - Valle d'Aosta	1	0.16	0.16	0.17	0.13	0.12	0.13
Lombardia	2	0.14	0.13	0.16	0.11	0.11	0.12
Trentino-Alto Adige	3	0.08	0.07	0.08	0.08	0.08	0.08
Veneto	4	0.15	0.14	0.16	0.12	0.11	0.13
Friuli Venezia Giulia	5	0.14	0.13	0.15	0.12	0.10	0.13
Liguria	6	0.14	0.12	0.15	0.11	0.10	0.12
Emilia Romagna	7	0.14	0.13	0.16	0.11	0.11	0.12
Toscana	8	0.15	0.14	0.16	0.12	0.11	0.13
Umbria	9	0.20	0.17	0.23	0.15	0.13	0.17
Marche	10	0.17	0.15	0.18	0.12	0.11	0.13
Lazio	11	0.19	0.17	0.22	0.13	0.12	0.15
Abruzzo	12	0.20	0.17	0.22	0.13	0.12	0.15
Molise	13	0.18	0.18	0.18	0.12	0.12	0.12
Campania	14	0.24	0.22	0.26	0.15	0.14	0.16
Puglia	15	0.24	0.21	0.26	0.16	0.14	0.17
Basilicata	16	0.20	0.16	0.23	0.13	0.11	0.14
Calabria	17	0.22	0.19	0.25	0.15	0.13	0.16
Sicilia	18	0.23	0.21	0.24	0.15	0.15	0.16
Sardegna	19	0.22	0.20	0.24	0.15	0.13	0.17
Italy		0.18	0.16	0.20	0.13	0.12	0.14

a strong correlation between the Head Count Ratio and the Gap. As the average level of dissatisfaction increases, the distance between the “unsatisfied” and the others becomes larger, giving evidence of a social polarization process, particularly affecting southern regions. Focusing on males and females separately reveals other interesting features in the territorial pattern of subjective well-being. Regional Head Count Ratios and Gaps are systematically higher for females than for males, revealing a kind of “gender polarization” across the country. Female subjective well-being differentiates regions more neatly than male scores at the extent that regions form quite separated clusters, enforcing the evidence of strong variations in the structure of subjective well-being along the North–South axis. Again, Trentino-Alto Adige is an exception, in that the difference between males and females is very small.

Remark on Missing Data As stated, records with missing data have been excluded by the analysis. Given the aim of the paper (to present the essentials of a new evaluation methodology), this is an acceptable choice. Indeed, an interesting feature of the methodology is that missing data could be handled quite easily. Each statistical unit in the population is assigned to an element of the profile poset and his/her well-being equals the value of the identification function on that element. When some components of a profile are missing, the statistical unit can only be associated to a subset of the profile poset, comprising the profiles compatible with the available information. Consequently, a range of possible well-being scores may be associated to the statistical unit. Similarly, a range of variation for the overall well-being score could be also derived. Due to the limited space available, here we cannot pursue this analysis further.

5 Conclusion

In this paper, we have introduced and applied to real data a new methodology for multidimensional evaluation with ordinal data, that overcomes the limitation of approaches based on counting or on scaling of ordinal variables. The proposed methodology exploits results from partial order theory and produces synthetic indicators with no variable aggregation.⁸ The approach is still under development, particularly to give it sound mathematical foundations, to tune it towards real applications and to overcome the computational limitations due to sampling from the set of linear extensions. Future and broader applications to real data will determine whether the methodology is valuable. The issue of well-being evaluation

⁸It is of interest to notice that in standard multivariate approaches, aggregation often exploits interdependencies among variables. Unfortunately, in quality-of-life studies, it turns out that interdependencies may be quite weak. Our approach, which is multidimensional in nature, overcomes this issue by addressing the evaluation problem as a problem of multidimensional comparison.

is gaining importance day by day, for both scholars and policy-makers. We hope to be contributing to address the problem in a more effective way.

References

1. Alkire, S., Foster, J.: Counting and multidimensional poverty measurement. *J. Public Econ.* **96**(7–8), 476–487 (2011)
2. Bublely, R., Dyer, M.: Faster random generation of linear extensions. *Discret. Math.* **201**, 81–88 (1999)
3. Fattore, M., Arcagni, A.: PARSEC: an R package for poset-based evaluation of multidimensional poverty. In: Bruggemann, R., Carlsen, L. (eds.) *Multi-Indicator Systems and Modelling in Partial Order*. Springer, New York (2014)
4. Fattore, M., Brueggemann, R., Owsiniński, J.: Using poset theory to compare fuzzy multidimensional material deprivation across regions. In: Ingrassia, S., Rocci, R., Vichi, M. (eds.) *New Perspectives in Statistical Modeling and Data Analysis*. Springer, New York (2011)
5. Fattore, M., Maggino, F., Greselin, F.: Socio-economic evaluation with ordinal variables: integrating counting and poset approaches. *Statistica Applicazioni* 31–42 (2011, Special Issue)
6. Fattore, M., Maggino, F., Colombo, E.: From composite indicators to partial order: evaluating socio-economic phenomena through ordinal data. In: Maggino, F., Nuvolati, G. (eds.) *Quality of Life in Italy: Research and Reflections, Social Indicators Research Series*, vol. 48. Springer, New York (2012)
7. Neggers, J., Kim, H.S.: *Basic Posets*. World Scientific, Singapore (1998)
8. R Core Team R: *A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna. ISBN 3-900051-07-0 (2013). <http://www.R-project.org/>.

Composite Indicator of Social Inclusion for the EU Countries

Francesca Giambona and Erasmo Vassallo

Abstract

Social inclusion is one of the key challenges of the European Union Sustainable Development Strategy (EU SDS). We use the main indicators identified by Eurostat within the operational objectives of the specific European policies to measure social inclusion for the 27 member countries of the European Union. In particular, we aggregate four basic indicators in a multiplicative composite indicator via a DEA-BoD approach with weights determined endogenously with proportion constraints. We obtain a score of social inclusion that allows us to grade the 27 EU countries from 2006 to 2010. In this way, we highlight the specific role played by the four indicators in determining improvements and deteriorations of social inclusion during the European phase of the financial and economic crisis.

1 Introduction

The European Commission defines social inclusion as “a process which ensures that those at risk of poverty and social exclusion gain the opportunities and resources necessary to participate fully in economic, social and cultural life and to enjoy a standard of living and well-being that is considered normal in the society in which they live. It ensures that they have a greater participation in decision making which affects their lives and access to their fundamental rights” [1, p. 10]. Social inclusion is, therefore, complementary concept to social exclusion, that is: “a process

F. Giambona (✉) • E. Vassallo

Dipartimento di Scienze Economiche, Aziendali e Statistiche, Università degli Studi di Palermo, Palermo, Italy

e-mail: fgiambona78@gmail.com; erasmo.vassallo@unipa.it

whereby certain individuals are pushed to the edge of society and prevented from participating fully by virtue of their poverty, or lack of basic competencies and lifelong learning opportunities, or as a result of discrimination. This distances them from job, income and education and training opportunities, as well as social and community networks and activities. They have little access to power and decision making bodies and thus feel powerless and unable to take control over the decisions that affect their day to day lives” [1, p. 10]. Multidimensionality of the social inclusion concept is evident. This complexity is manifested in the use of multiple measures. For example, the Lisbon European Council already in 2000 has suggested the use of a set of indicators to measure progress in relation to poverty and social exclusion [2]. The theme and its measurement is highly topical; it represents one of the key challenges of the European Union Sustainable Development Strategy (EU SDS) in order to achieve a socially inclusive society, to increase the quality of life of citizens and individual well-being. In particular, reduction of poverty and reduction of social exclusion is one of the five main targets of “Europe 2020”, that is, the EU’s strategy for a sustainable and inclusive growth for the next years [3].

In line with this strategy, the European debate has focused on four main pillars or dimensions from which derive four specific policy actions for a higher social inclusion of individuals: (a) reduction of monetary poverty; (b) improvement of living conditions; (c) greater access to labour markets and (d) better education [4–7]. However, in order to monitor levels, deteriorations and improvements of social inclusion in Europe, it appears useful to identify a single measure that can summarize these four pillars to determine a ranking of countries and, consequently, to compare the intensity of social inclusion. For this purpose, we construct a composite indicator of social inclusion at macro level for the 27 member countries of the European Union by aggregating the pillars represented, each, by one specific basic indicator related to the operational objectives and targets (so-called level 2) of the European policies identified by Eurostat [6, 8]. These four indicators are listed below.

- 1(pove). People at risk of poverty after social transfers (percentage of total population). This indicator measures the share of persons at risk of monetary poverty. Persons are at risk of poverty if their equivalized disposable income is below the risk-of-poverty threshold, which is set at 60 % of the national median after social transfers.
- 2(depr). Severely materially deprived people (percentage of total population). It covers issues relating to economic strain and durables. Severely materially deprived persons have living conditions greatly constrained by a lack of resources and cannot afford at least four of the following: to pay rent or utility bills; to keep their home adequately warm; to pay unexpected expenses; to eat meat, fish or a protein equivalent every second day; a week holiday away from home; a car; a washing machine; a colour TV or a telephone.

- 3(work). People living in households with very low work intensity (percentage of total population). Persons are defined as living in households with very low work intensity if they are aged 0–59 and the working age members in the household worked less than 20 % of their potential during the past year.
- 4(educ). Early school leavers (percentage of total). It is defined as the percentage of the population aged 18–24 with at most lower secondary education and not in further education or training.¹

Synthesis of the four indicators in a single measurement of social inclusion appears of great relevance to provide a unique information about the status of the country and to identify directions for improvement in the light of the European policies. Anyhow, aggregation of indicators is not a trivial issue and, still, the best choice of the weights is a topic of interest. Literature proposes two main type of aggregation, additive and multiplicative, and several weighting methods such as equal weights, weights based on statistical models and weights based on public/expert opinion [10]. For our goals, a multiplicative DEA-like model in a benefit-of-doubt (BoD) approach seems the most appropriate alternative. In effect, this method allows varying weights determined endogenously according to an optimal solution looking for the best possible outcome for the country under analysis. In our case, this implies that the composite indicator of social inclusion combines the four sub-indicators in the best interest of the country. Of course, this does not mean that the obtained weighting, certainly optimal from a mathematical point of view, is also optimal from the point of view of politics, but it is certainly independent from subjective experiences of the decision maker.

Finally, the composite indicator of social inclusion is calculated over the years 2006–2010, the longest available period at the time of writing without missing data or data breaks for all the 27 EU countries, precisely: Austria (AT), Belgium (BE), Bulgaria (BG), Cyprus (CY), Czech Republic (CZ), Germany (DE), Denmark (DK), Estonia (EE), Greece (EL), Spain (ES), Finland (FI), France (FR), Hungary (HU), Ireland (IE), Italy (IT), Lithuania (LT), Luxembourg (LU), Latvia (LV), Malta (MT), Netherlands (NL), Poland (PL), Portugal (PT), Romania (RO), Sweden (SE), Slovenia (SI), Slovakia (SK) and United Kingdom (UK). Interestingly, this period covers the severe phase of the financial and economic crisis in Europe.

¹We note that although the multidimensional nature of social inclusion is widely shared, some dimensions are neglected in the time of measurement. In effect, EU social indicators are much better developed for poverty, material deprivation, labour market and level of education than for political or cultural dimensions [9]. However, the use of the same four basic indicators selected by Eurostat allows us to maintain a strong connection with the objectives set by the European strategies and allows us an easy comparability of the results.

2 Multiplicative Composite Indicator

The construction of a composite indicator is not a trivial argument. The advantage of an immediate synthesis presents also some disadvantages and risks. Synthesis, inevitably, reduces information and, often, uses subjective options open to criticism. However, in many cases, it may represent a very useful tool to provide a clear view of analysis for defining appropriate policies.² The weighted arithmetic mean

$\sum_{j=1}^n w_j I_{ij}$ is the easiest way to create a composite indicator, though the traditional

additive approach shows some important limitations. An undesirable feature of additive aggregations is the full compensability that they imply, since low values in some dimensions can be compensated by high values in other dimensions. On the contrary, the geometric aggregation is a less compensatory approach, and it offers better incentives to countries to improve underperforming dimensions. In fact, marginal utility on the multiplicative composite indicator is higher if the indicator value is low, this reflecting the auspices of the European policy for a progressive improvement of all the dimensions of social inclusion [10].³ With this primary motivation, we use a weighted product method in construction of the composite indicator of social inclusion (SI), formulated as

$$SI_i = \prod_{j=1}^n I_{ij}^{w_j} \quad (1)$$

where I_{ij} is the j -th basic indicator of social inclusion ($j = 1, \dots, n$) for the i -th country ($i = 1, \dots, m$) with weight w_j . Here, $n = 4$ and $m = 27$. SI_i looks like a geometric mean.⁴

²Many composite indicators exist in literature with varying degrees of methodological complexity. For example, the “Corruption Perceptions Index” by Transparency International [11], the “Human Development Index” by UN [12], the “Composite Leading Indicators” by OECD [13] and so on. A good starting point on the issue is OECD [10].

³Geometric aggregation is a good compromise between methods with full compensability and non-compensatory approaches, for example, MCDA (Multiple-Criteria Decision Analysis) [14]. In general terms, geometric aggregation is preferable to approaches MCDA because it can lead to the minimum information loss [15]. Furthermore, using the social choice theory, Ebert and Welsch [16] found that geometric aggregation is particularly relevant in composite indicator construction when the ordinal information is involved.

⁴With multiplicative aggregation, the sub-indicators must be larger than 1 otherwise the logarithmic transformation obtains negative values. Therefore, it is necessary to normalize the data of the four basic indicators extracted from Eurostat. In particular, we use a min–max transformation in a continuous scale from 2 (minimum) to 10 (maximum) where higher values correspond to better social inclusion. In other words, we apply the transformation: $(\max(y) - y) / (\max(y) - \min(y)) \cdot 8 + 2$. Furthermore, we note that in this way the direction of the sub-indicators is reversed so that higher values represent, in more intuitive terms, greater social inclusion and not greater social exclusion as in the original scale of the Eurostat indicators. For our purposes, given the techniques used here, this does not distort the final result as it will be clear later. Finally, it should be noted

Anyhow, results can depend strongly on the selected weights. To avoid subjective choices easily criticized, in this paper the weights w_j are determined endogenously with an automatic mechanism based on a multiplicative optimization model similar to a DEA-BoD model written as follows:

$$\begin{aligned}
 SI_i &= \max_w \prod_{j=1}^n I_{ij}^{w_j} \\
 \text{s.t.} & \\
 \prod_{j=1}^n I_{ij}^{w_j} &\leq e \\
 w_j &\geq 0
 \end{aligned} \tag{2}$$

where e is the Napier's constant [17–19].⁵

In this way, the composite indicator is obtained by multiplying the four basic indicators of social inclusion with weights calculated in the best possible conditions, i.e. increasing as much as possible the composite score of social inclusion for a given country.⁶ In short, a low value of the composite indicator SI_i , then low social inclusion for the i -th country, is due to low values of the indicators that compose it and not attributable to specific weights, always calculated to obtain the best, i.e. maximum, possible result for the i -th country compared to the benchmark country. In fact, the composite indicator is defined as the ratio of a country's actual performance to its benchmark performance and its maximum score is at most equal to $e = 2.718281828$.⁷

Searching for the best values, the optimization problem could give zero weight to some indicators and attribute too much weight to other indicators (curse of dimensionality): this is not desirable if all the four dimensions are theoretically relevant. Besides, this involves no unique ranking and, in some special cases, no

immediately that correlations among the four sub-indicators are moderate; in fact, there is no risk of double-counting: this is an ideal condition in the construction of a composite indicator.

⁵The Benefit-of-the-Doubt (BoD) logic assumes a favourable judgement in the absence of full evidence using a model similar to Data Envelopment Analysis (DEA) [20]. In fact, we are not sure about the appropriate weights, rather we look for BoD weights such that the country's performance of social inclusion is as high as possible [21, 22, 18, 17, 23, 24, 25, 26, 27, 28]. In brief, model (2) is like an output-oriented DEA model where indicators y are outputs and a variable always equal to one is the only input: it is the Koopmans "helmsman", by which countries have an apparatus responsible for the conduct of their social policies [29]. Therefore, the social inclusion performance is evaluated in terms of the ability of the helmsman in each country to maximize the levels of the four basic indicators (obviously normalized in terms of social inclusion) [30].

⁶It should be noted that DEA typically does not require normalization of the data, made here for the unique needs of greater convenience of analysis. It is not even mandatory that the unit of measurement is identical, since the weights take into account the unit of measurement of the sub-indicators [30].

⁷This means a preference for an internal benchmark as the best practice country, rather than an external benchmark that could not be realistically achievable in specific local contexts.

feasible solution. For these reasons, we add specific constraints on the weights; in particular, we add proportion constraints to the model (2):

$$\left(\prod_{j=1}^n I_{ij}^{w_j} \right)^L \leq I_{ij}^{w_j} \leq \left(\prod_{j=1}^n I_{ij}^{w_j} \right)^U \quad (3)$$

where L and U , ranging between 0 and 1, represent the lower and upper bound for the contribution (in percentage terms) of the j -th sub-indicator [17].⁸ To avoid zero weights, we impose $L = 20\%$ with U determined accordingly for all the 27 countries and the four sub-indicators.⁹ Different limits do not give solution to the mathematical problem or lead to an excessive number of benchmark countries or, also, to an unjustifiable imbalance of the indicators' role.

We note that the multiplicative optimization problem (2) and (3) is nonlinear, then we solve the equivalent linear problem by taking logarithms with base e ; at the end, it is easy to obtain the original multiplicative indicator SI . If social inclusion is higher, SI value is higher, where $SI = e$ indicates the benchmark country. We remark that the maximum score is always e , but the empirical minimum can change; so, for easy and accurate comparison, the scores are normalized by their range of variation.

3 Results

Table 1 shows the social inclusion (SI) score for the 27 EU countries over the years 2006–2010. For reasons of easy comparability, SI is normalized by its range; therefore, the score is now between 0 (the lowest level of social inclusion) and 1 (the maximum level of social inclusion, i.e. benchmark country). For convenience of the reader, here is also shown the additive case and the simple arithmetic mean without weights, but for the sake of brevity we focus only on the multiplicative scores. In detail, Table 2 shows the final multiplicative scores and their decomposition in the four sub-indicators of social inclusion just at the beginning and at the end of the period.

⁸Without constraints on the weights, the DEA model could reset the contribution of the underperforming dimensions to find the best solution. Thus, the results could depend even on a single indicator and, consequently, we could have a large number of insignificant benchmarks. This event occurs more likely when the sample is not large as in our case. Specific constraints on the weights and use of a few indicators, compared to the number of countries, avoid this “curse of dimensionality” [20]. Here, in particular, we use proportion constraints which offer a very intuitive interpretation that we consider preferable to the available alternatives such as absolute, relative or ordinal restrictions.

⁹For example, if in a given country three indicators reach the minimum contribution of 20%, the contribution of the fourth indicator will necessarily have an upper bound of 40%. In fact, the specification of the upper bound is not necessary since the sum of the contributions of the four indicators must be 100%.

Table 1 Normalized social inclusion scores for the EU countries (2006–2010)

Countries	Multiplicative aggregation				Additive aggregation				Arithmetic mean									
	2006	2007	2008	2009	2010	+/-	2006	2007	2008	2009	2010	+/-	2006	2007	2008	2009	2010	+/-
AT: Austria	0.944	0.941	0.890	0.921	0.928	-	0.925	0.919	0.854	0.893	0.898	-	0.893	0.899	0.800	0.827	0.830	-
BE: Belgium	0.423	0.632	0.577	0.739	0.750	+	0.388	0.565	0.490	0.665	0.671	+	0.544	0.611	0.500	0.549	0.543	-
BG: Bulgaria	0.000	0.000	0.000	0.000	0.000	-	0.000	0.000	0.000	0.000	0.000	-	0.000	0.000	0.000	0.000	0.020	+
CY: Cyprus	0.955	0.943	0.913	0.874	0.791	-	0.921	0.915	0.880	0.825	0.718	-	0.844	0.861	0.808	0.732	0.636	-
CZ: Czech Republic	0.970	0.991	1.000	1.000	1.000	+	0.964	0.988	1.000	1.000	1.000	+	0.944	0.976	1.000	1.000	1.000	+
DE: Germany	0.586	0.772	0.584	0.759	0.756	+	0.537	0.709	0.497	0.687	0.676	+	0.631	0.687	0.494	0.567	0.561	-
DK: Denmark	0.921	0.880	0.897	0.883	0.862	-	0.899	0.842	0.866	0.842	0.809	-	0.880	0.826	0.809	0.744	0.703	-
EE: Estonia	0.809	0.782	0.775	0.694	0.731	-	0.748	0.714	0.705	0.604	0.648	-	0.707	0.717	0.677	0.553	0.554	-
EL: Greece	0.671	0.686	0.622	0.645	0.423	-	0.585	0.601	0.533	0.548	0.335	-	0.553	0.597	0.454	0.456	0.365	-
ES: Spain	0.617	0.590	0.466	0.412	0.242	-	0.514	0.491	0.355	0.303	0.169	-	0.526	0.540	0.380	0.276	0.178	-
FI: Finland	0.919	0.921	0.919	0.883	0.889	-	0.894	0.894	0.890	0.840	0.844	-	0.862	0.871	0.831	0.749	0.749	-
FR: France	0.877	0.868	0.838	0.850	0.825	-	0.839	0.827	0.789	0.800	0.764	-	0.794	0.798	0.728	0.702	0.652	-
HU: Hungary	0.548	0.766	0.493	0.662	0.662	+	0.473	0.702	0.412	0.576	0.575	+	0.455	0.674	0.404	0.458	0.460	+
IE: Ireland	0.587	0.574	0.138	0.121	0.173	-	0.516	0.505	0.160	0.147	0.165	-	0.515	0.562	0.385	0.283	0.216	-
IT: Italy	0.620	0.651	0.571	0.636	0.570	-	0.530	0.564	0.485	0.536	0.462	-	0.449	0.520	0.338	0.398	0.358	-
LT: Lithuania	0.677	0.799	0.755	0.611	0.371	-	0.591	0.733	0.693	0.513	0.285	-	0.541	0.723	0.681	0.458	0.294	-
LU: Luxembourg	0.989	0.989	1.000	0.924	0.934	-	0.974	0.982	1.000	0.892	0.900	-	0.926	0.954	0.957	0.831	0.863	-
LV: Latvia	0.478	0.640	0.119	0.063	0.074	-	0.393	0.549	0.121	0.060	0.048	-	0.389	0.539	0.301	0.150	0.000	-
MT: Malta	0.269	0.291	0.110	0.096	0.178	-	0.249	0.273	0.117	0.100	0.160	-	0.453	0.484	0.299	0.253	0.262	-
NL: Netherlands	0.878	0.927	0.944	0.932	0.970	+	0.848	0.905	0.925	0.907	0.954	+	0.846	0.900	0.883	0.826	0.888	+
PL: Poland	0.567	0.762	0.754	0.768	0.673	+	0.484	0.697	0.684	0.697	0.581	+	0.431	0.644	0.593	0.615	0.546	+
PT: Portugal	0.313	0.373	0.299	0.418	0.457	+	0.261	0.311	0.220	0.310	0.340	+	0.398	0.443	0.265	0.258	0.241	-
RO: Romania	0.208	0.227	0.207	0.291	0.082	-	0.183	0.198	0.122	0.186	0.049	-	0.210	0.277	0.021	0.047	0.022	-

(continued)

Table 1 (continued)

Countries	Multiplicative aggregation					Additive aggregation					Arithmetic mean							
	2006	2007	2008	2009	2010	+/-	2006	2007	2008	2009	2010	+/-	2006	2007	2008	2009	2010	+/-
SE: Sweden	0.971	1.000	1.000	0.931	0.953	-	0.957	1.000	1.000	0.902	0.927	-	0.927	0.995	0.966	0.828	0.865	-
SI: Slovenia	1.000	1.000	0.976	0.974	0.925	-	1.000	1.000	0.962	0.963	0.894	-	1.000	1.000	0.935	0.935	0.853	-
SK: Slovakia	0.971	0.987	1.000	0.948	0.890	-	0.953	0.978	1.000	0.930	0.852	-	0.919	0.963	0.981	0.889	0.785	-
UK: United Kingdom	0.638	0.716	0.584	0.651	0.639	+	0.564	0.638	0.499	0.558	0.538	-	0.543	0.595	0.385	0.411	0.409	-
Mean	0.682	0.730	0.646	0.655	0.620	-	0.637	0.685	0.602	0.603	0.565	-	0.636	0.691	0.588	0.548	0.513	-
Median	0.671	0.772	0.754	0.739	0.731	+	0.585	0.709	0.684	0.665	0.648	+	0.553	0.687	0.593	0.553	0.546	-
Std. deviation	0.269	0.254	0.315	0.301	0.312	+	0.278	0.262	0.320	0.301	0.313	+	0.249	0.233	0.288	0.273	0.288	+
Coeff. of variation	0.395	0.348	0.487	0.460	0.503	+	0.437	0.383	0.532	0.500	0.554	+	0.392	0.337	0.490	0.499	0.561	+

Table 2 Multiplicative social inclusion scores and decomposition for the EU countries (2006 and 2010)

Countries	Multiplicative scores (year 2006)					Multiplicative scores (year 2010)				
	1. Povc	2. Depr	3. Work	4. Educ	Score	1. Povc	2. Depr	3. Work	4. Educ	Score
AT: Austria	1.216	1.479	1.216	1.216	2.659	1.215	1.476	1.215	1.215	2.649
BE: Belgium	1.161	1.273	1.161	1.230	2.110	1.199	1.437	1.199	1.199	2.476
BG: Bulgaria	1.107	1.107	1.107	1.226	1.664	1.118	1.118	1.250	1.118	1.748
CY: Cyprus	1.217	1.217	1.481	1.217	2.670	1.203	1.203	1.446	1.203	2.516
CZ: Czech Republic	1.485	1.219	1.219	1.219	2.687	1.221	1.221	1.221	1.492	2.718
DE: Germany	1.208	1.358	1.179	1.179	2.282	1.199	1.438	1.199	1.199	2.482
DK: Denmark	1.214	1.473	1.214	1.214	2.635	1.209	1.462	1.209	1.209	2.585
EE: Estonia	1.203	1.203	1.447	1.203	2.517	1.197	1.433	1.197	1.197	2.457
EL: Greece	1.188	1.412	1.188	1.188	2.371	1.166	1.166	1.360	1.166	2.158
ES: Spain	1.183	1.183	1.399	1.183	2.314	1.147	1.315	1.147	1.147	1.983
FI: Finland	1.214	1.473	1.214	1.214	2.633	1.212	1.468	1.212	1.212	2.610
FR: France	1.210	1.463	1.210	1.210	2.589	1.206	1.454	1.206	1.206	2.548
HU: Hungary	1.175	1.175	1.175	1.381	2.242	1.190	1.190	1.190	1.417	2.391
IE: Ireland	1.179	1.391	1.179	1.179	2.283	1.139	1.297	1.139	1.139	1.916
IT: Italy	1.183	1.400	1.183	1.183	2.318	1.181	1.396	1.181	1.181	2.301
LT: Lithuania	1.189	1.189	1.189	1.414	2.377	1.161	1.161	1.161	1.348	2.108
LU: Luxembourg	1.220	1.220	1.489	1.220	2.707	1.216	1.478	1.216	1.216	2.655
LV: Latvia	1.167	1.167	1.363	1.167	2.168	1.127	1.127	1.127	1.271	1.820
MT: Malta	1.143	1.306	1.143	1.143	1.948	1.139	1.298	1.139	1.139	1.920
NL: Netherlands	1.463	1.210	1.210	1.210	2.590	1.219	1.485	1.219	1.219	2.689

(continued)

Table 2 (continued)

Countries	Multiplicative scores (year 2006)				Multiplicative scores (year 2010)				Score	
	1. Pove	2. Depr	3. Work	4. Educ	1. Pove	2. Depr	3. Work	4. Educ		
PL: Poland	1.177	1.177	1.177	1.386	1.191	1.191	1.191	1.420	2.262	2.401
PT: Portugal	1.148	1.148	1.318	1.148	1.170	1.369	1.170	1.170	1.994	2.191
RO: Romania	1.135	1.135	1.288	1.135	1.128	1.128	1.273	1.128	1.883	1.828
SE: Sweden	1.219	1.485	1.219	1.219	1.217	1.482	1.217	1.217	2.688	2.673
SI: Slovenia	1.221	1.221	1.221	1.492	1.215	1.476	1.215	1.215	2.718	2.645
SK: Slovakia	1.222	1.219	1.480	1.219	1.212	1.212	1.212	1.468	2.687	2.612
UK: United Kingdom	1.185	1.404	1.185	1.185	1.188	1.412	1.188	1.188	2.337	2.368

In general terms, between 2006 and 2010, only a few countries appear closer to their benchmark of social inclusion (see, sign + in Table 1). This is enough to raise median of *SI* from 0.671 in 2006 to 0.731 in 2010. However, also variability (coefficient of variation) increases from 0.395 to 0.503, indicating a higher dispersion of the scores with a relative worsening in some countries (mean is reduced from 0.682 to 0.620). Bulgaria occupies the last position in all the years due to poor performances of the four dimensions of social inclusion. Conversely, in 2006, we have Slovenia in the first position (low levels of poverty and school leavers) whereas, in 2007, the leadership is divided between Slovenia and Sweden (especially for their low levels of poverty).

In 2008, the best countries are three: Czech Republic, Sweden and Slovakia (once again, the low levels of poverty and school leavers play an important role). From this moment on, Czech Republic will be the only benchmark in 2009 and 2010. From 2006 to 2010, only 8 countries are closer to the benchmark, in particular, Belgium (from 0.423 to 0.750, with a strong improvement of the materially deprived people). Differently, it is interesting to note a substantial (relative) worsening in Ireland (from 0.587 to 0.173), Latvia (from 0.478 to 0.074), Spain (from 0.617 to 0.242), Lithuania (from 0.677 to 0.371) and Greece (from 0.671 to 0.423). In short, Lithuania loses 8 positions (from 13 to 21), Greece and Ireland lose 6 positions (from 14 to 20 and from 18 to 24, respectively), whereas Belgium gains 10 positions (from 23 to 13) and the Netherlands gains 8 positions (from 10 to 2). Italy loses 2 positions (from 16 in 2006 to 18 in 2010).

Further consideration should be made on the consistency of the composite indicator with respect to the information provided in the so-called headline indicator (*HI*), that is, “people at risk of poverty or social exclusion”, used by Eurostat to represent in a simple and single measure the level of social inclusion.

The values are consistent but sufficiently different. In effect, the headline indicator is not a composite indicator; in addition, *SI* includes some aspects of education. In fact, *HI* considers people who are at risk of poverty *or* severely materially deprived *or* living in households with very low work intensity and, in case of intersections, a person is counted only once; so, one can say that *SI* is a more broad measure and, at the same time, more fitting for a given country than the headline indicator.

At a glance, the combination of the four sub-indicators of social inclusion is not trivial denoting specificities and characteristics of each country that *SI* considers properly through differentiated weights among indicators and countries without applying external information. Of course, this does not imply that the achieved results are preferable or better from the point of view of politics. Moreover, the

decomposition of the *SI* score is useful to highlight, country by country, the most critical aspects of social inclusion to which policies should be directed.¹⁰

4 Conclusion

The European Union has adopted the “Europe 2020” strategy aimed at a sustainable and inclusive growth in the 27 member countries, in order to achieve high levels of employment, productivity and social cohesion through five objectives to be reached by 2020, including reduction of the school drop-out rates and poverty to achieve higher levels of social inclusion that, moreover, is a key target of the EU SDS [31]. In fact, social inclusion is a complex concept with a multidimensional perspective that, for purposes of policy, makes useful the synthesis in a single measure.

Therefore, in this paper, we have focused on the construction of a composite indicator of social inclusion able to represent the position of each country and to identify directions for improvement with respect to some benchmarks. In particular, to stay close enough to the strategy of the EU and its translation in the indicators proposed by Eurostat, also to promote greater comparability of our results, we have used data corresponding to the four main operational objectives (or “level 2” targets) as defined in the Sustainable Development Strategy, and related to poverty, material deprivation, labour market and education, although other political or cultural dimensions remain neglected [32]. Our composite indicator of social inclusion (*SI*) is obtained by aggregating these four dimensions using appropriate weights determined with an automatic procedure based on a multiplicative DEA-BoD method with proportion constraints. This optimizes mathematically the results for every country without resorting to external judgments of experts, even if the outcome may not be politically satisfying. Furthermore, the multiplicative aggregation emphasizes the improvements of the countries with worse conditions of social inclusion; in other words, this aggregation implies a partial compensability of the four sub-indicators offering better motivations to improve underperforming dimensions.

Data are collected from 2006 to 2010 for all the 27 EU member countries in the largest time interval available with no missing data at the time of writing. In general terms, from 2006 to 2010, the results show a small shift toward the benchmark values of social inclusion, even if there are significant deteriorations

¹⁰It is appropriate to make a final comment about the robustness of the results. Some countries may be outliers and strongly influence the score of *SI*. To verify this vulnerability of the results, we have repeated $m = 27$ times the calculation of *SI* removing each time a different country. The impact of the j -th missing country was measured through the sum of the $m - 1$ squared differences between the score of the i -th countries ($i \neq j$) computed including the j -th country and that one computed excluding the j -th country. So, we obtain 27 values each representing the influence on the *SI* score of the country from time to time excluded from the calculation. The differences are very small in many cases and, sometimes, completely negligible, even when they involve the benchmark countries.

in some countries and greater variability of the SI scores among all the 27 European countries. In fact, an important role is played by changes in the levels of monetary poverty, here measured in relative (not absolute) terms as required by the European policy. Finally, we note that some of these 27 countries worsen their levels of social inclusion since 2008 in relation to the European contagion of the financial and economic crisis.

References

1. European Commission: Joint Report on Social Inclusion. Directorate General for Employment, Social Affairs and Equal Opportunities. Bruxelles (2004)
2. European Council: Lisbon European Council 23 and 24 march 2000. Presidency Conclusions, Bruxelles (2000)
3. European Commission: Europe 2020. A Strategy for Smart, Sustainable and Inclusive Growth. Communication COM (2010). Bruxelles (2010a)
4. Atkinson, A.B., Cantillon, B., Marlier, E., Nolan, B.: Social Indicators: The EU and Social Inclusion. Oxford University Press, Oxford (2002)
5. Atkinson, A.B., Marlier, E., Nolan, B.: Indicators and targets for social inclusion in the European Union. *J. Common Mark. Stud.* **42**(1), 47–75 (2004)
6. Eurostat: Sustainable Development in the European Union. Eurostat, Bruxelles (2011a)
7. Marlier, E.: Setting targets: the use of indicators. In: EAPN Network News, vol. 98, pp. 4–6 (2003)
8. Eurostat: Sustainable Development Indicators: Social Inclusion. Eurostat, Bruxelles (2011b)
9. Berger-Schmitt, R., Noll, H.H.: Conceptual framework and structure of a European system of social indicators. In: EuReporting Working Paper, no. 9, Mannheim, ZUMA (2000)
10. OECD: Handbook on Constructing Composite Indicators: Methodology and User Guide. OECD, Paris (2008)
11. Transparency International: Annual Report. Transparency International, New York (2011)
12. UN: Human Development Report. UNDP, New York (2011)
13. OECD: Composite Leading Indicators. OECD, Paris (2012)
14. Munda, G., Nardo, M.: On the methodological foundations of composite indicators used for ranking countries. In: OECD/JRC Workshop on composite indicators of country performance, Ispra, Italy, 12 May 2003
15. Zhou, P., Ang, B.W.: Comparing MCDA aggregation methods in constructing composite indicators using the Shannon-Spearman measure. *Soc. Indic. Res.* **94**, 83–96 (2009)
16. Ebert, U., Welsch, H.: Meaningful environmental indices: a social choice approach. *J. Environ. Econ. Manag.* **47**, 270–283 (2004)
17. Cherchye, L., Moesen, W., Rogge, N., van Puyenbroec, T., Saisana, M., Saltelli, A., Liska, R., Tarantola, S.: Creating composite indicators with DEA and robustness analysis: the case of the technology achievement index. *J. Oper. Res. Soc.* **59**, 239–251 (2008)
18. Zhou, P., Ang, B.W., Poh, K.L.: A mathematical programming approach to constructing composite indicators. *Ecol. Econ.* **62**, 291–297 (2007)
19. Zhou, P., Fan, L., Zhou, D.: Data aggregation in constructing composite indicators: a perspective of information loss. *Expert Syst. Appl.* **37**, 360–365 (2010)
20. Coelli, T.J., Rao, D.P., O'Donnell, C.J., Battese, G.E.: An Introduction to Efficiency and Productivity Analysis. Springer, New York (2005)
21. Melyn, W., Moesen, W.: Towards a synthetic indicator of macroeconomic performance: unequal weighting when limited information is available. In: Public Economics Research Paper CES, vol. 17, Ces, KU Leuven (1991)
22. Despotis, D.K.: Measuring human development via data envelopment analysis: the case of Asia and the Pacific. *OMEGA Int. J. Manag. Sci.* **33**, 385–390 (2005)

23. Cherchye, L., Knox Lovell, C.A., Moesen, W., Van Puyenbroeck, T.: One market, one number? A composite indicator assessment of EU internal market dynamics. *Eur. Econ. Rev.* **51**, 749–779 (2007)
24. Cherchye, L., Moesen, W., Rogge, N., Van Puyenbroeck, T.: An introduction to ‘benefit of the doubt’ composite indicators. *Soc. Indic. Res.* **82**, 111–145 (2007)
25. Mohamad, N.: A linear programming formulation of macroeconomic performance: the case of Asia Pacific journal. *Matematika* **23**, 29–40 (2007)
26. Murias, P., de Miguel, J., Rodríguez, D.: A composite indicator for university quality assessment: the case of Spanish higher education system. *Soc. Indic. Res.* **89**, 129–146 (2008)
27. Zhou, P., Ang, B.W.: Linear programming models for measuring economy-wide energy efficiency performance. *Energy Policy* **36**, 2901–2906 (2008)
28. Zhou, P., Ang, B.W., Poh, K.L.: A survey of data envelopment analysis in energy and environmental studies. *Eur. J. Oper. Res.* **189**, 1–18 (2008)
29. Koopmans T.C. (ed.): *Activity Analysis of Production and Allocation*. Cowles Commission for Research in Economics, vol. 13. Wiley, New York (1951)
30. Knox Lovell, C.A., Pastor, J.T., Turner, J.A.: Measuring macroeconomic performance in the OECD: a comparison of European and Non-European countries. *Eur. J. Oper. Res.* **87**, 507–518 (1995)
31. European Council: *Conclusions*. Presidency Conclusions, Bruxelles (2010)
32. European Commission: *Joint Report on Social Protection and Social Inclusion*. Directorate General for Employment, Social Affairs and Equal Opportunities. Bruxelles (2010b)

A Well-Being Index Based on the Weighted Product Method

Matteo Mazziotta and Adriano Pareto

Abstract

It has long been accepted that the GDP per capita cannot alone explain the well-being in a geographical area. Several have been the attempts to construct alternative, non-monetary, indices of well-being by aggregating a variety of individual indicators that represent different dimensions of well-being. The most famous, in Italy, are the Index of Regional Quality of Development and the “Il Sole 24 Ore” Quality of Life Index. An issue often not solved, from a methodological point of view, concerns the comparability of the data over time. In this work, we propose a ‘static’ and a ‘dynamic’ well-being measure based on the application of the Jevons index to the socio-economic indicators. The obtained indices are closely related and allow synthetic spatial and temporal comparisons of the level of well-being.

1 Introduction

In the last 2 decades, many have been the attempts of different institutions (universities, statistics offices, international organizations) to construct composite indices of well-being, sustainable development or societal progress [1, 2]. In Italy, two interesting examples of this type are the Index of Regional Quality of Development (QUARS) proposed by the campaign “Sbilanciamoci!” and the Quality of Life Index published by the economic newspaper “Il Sole 24 Ore.”

Recently, the Italian National Institute of Statistics (Istat) has launched a series of studies for measuring equitable and sustainable well-being in Italy [6]. The aim

M. Mazziotta (✉) • A. Pareto

Italian National Institute of Statistics, Viale Oceano Pacifico 171, Rome, Italy

e-mail: mazziott@istat.it; pareto@istat.it

of the project, called BES (Benessere Equo Sostenibile), is to construct a set of measures of the various dimensions of well-being at regional level and for particular categories of people (e.g., male and female).

One of the main problems in order to construct composite indices is the choice of a method that allows comparisons over time. As is known, there are several procedures for the normalization of the data, most of which use ‘relative’ parameters (e.g., the average value, the minimum or the maximum of a given year). These parameters affect both the QUARS, which is based on z -scores, and the index proposed by “Il Sole 24 Ore,” that uses a function of ‘distance from the best performer.’ In the case of the Human Development Index (HDI), the problem has been overcome by using a re-scaling of the indicators in the range (0, 1) with limits independent from the observed values in a given year. This solution may lead to future values outside the range and the only alternative is to recalculate the index values for the past years [9].

In this paper, we propose an application of the Jevons index to the indicators of well-being that allows to build, for each unit, both a ‘static’ index, for regional comparisons, and a ‘dynamic’ index, for temporal comparisons, in a not full compensatory perspective. In Sect. 2, a description of the method is reported and in Sect. 3 an application to real data is proposed.

2 ‘Static’ and ‘Dynamic’ Well-Being Index

The Weighted Product (WP) method is one of the major techniques in composite index construction since it represents a trade-off solution between additive methods with full compensability and non-compensatory approaches [8]. Zhou et al. [11, 12] showed that the WP method may lead to a lower information loss in composite index construction compared to other aggregation methods.

When an unweighted geometric mean of ratios, such as the Jevons index, is computed, the obtained result satisfies many desirable properties from an axiomatic point of view [5].

Let us consider a set of individual indicators positively related with the well-being and let x_{ij}^t denote the value of the indicator j for the region i at time t , where $x_{ij}^t > 0$ ($j = 1, \dots, m$; $i = 1, \dots, n$; $t = t_0, t_1$). A ‘static’ well-being index may be defined as follows:

$$SWI_i^t = \prod_{j=1}^m \left(\frac{x_{ij}^t}{x_{rj}^t} 100 \right)^{\frac{1}{m}}$$

where x_{rj}^t is the reference or base value, e.g., the national average. Therefore, values of SWI that are higher (lower) than 100 indicate regions with above (below) average performance.

In order to compare the data from time t_0 to t_1 , for each region, we can construct a ‘dynamic’ well-being index given by:

$$DWI_i^{t_1/t_0} = \prod_{j=1}^m \left(\frac{x_{ij}^{t_1}}{x_{ij}^{t_0}} 100 \right)^{\frac{1}{m}}.$$

For the ‘circularity’ or ‘transitivity’ property of the index number theory, SWI and DWI are linked by the relation:

$$DWI_i^{t_1/t_0} = \left(\frac{SWI_i^{t_1}}{SWI_i^{t_0}} \right) DWI_r^{t_1/t_0}.$$

Note that the ‘dynamic’ well-being index is similar to the Canadian Index of Well-Being [7], except for the aggregation function. The Canadian approach is full compensatory since a simple arithmetic mean of ratios is used. We think that a multiplicative approach, such as in the new HDI [10], is preferable from both an axiomatic point of view (properties of the index) and a conceptual point of view (full compensability is not realistic) [4].

SWI and DWI are meaningful only for positive indicators. They give more weight to the small values and implicitly penalize the ‘unbalance’ among components [3].

3 An Application to the Italian Regions

In order to show the calculation of SWI and DWI, we consider a set of indicators of well-being in the Italian regions in 2005 and 2009.

The variables used are: Sporting activities, Close to supermarkets, Green space, Public transport, Parking provision, Children’s services, Elderly home care.

The data matrix is reported in Table 1.

Table 2 presents the results. Note that the base value of the ‘static’ indices (SWI_{05} and SWI_{09}), for each region, is the national average (Italy), while the base of the ‘dynamic’ index ($DWI_{09/05}$) is the value for the year 2005.

Moreover, the value of the ‘dynamic’ well-being index for Italy ($DWI_{09/05} = 108.5$) allows to obtain all the other ‘dynamic’ indices on the base of the ‘static’ ones. For example, we have the following result for Toscana:

$$DWI_{09/05} = \left(\frac{107.3}{113.4} \right) 108.5 = 102.6.$$

As we can see from the table, not necessarily each relative increase corresponds to an absolute one and vice versa. Indeed, from 2005 to 2009, Toscana shows a reduction of the level of well-being compared to the national average ($SWI_{09} - SWI_{05} = -6.2$), though the values of the individual indicators, on the whole, are increased ($DWI_{09/05} = 102.6$). This is due to a greater rise of the

Table 1 Individual indicators of well-being in the Italian regions (years 2005, 2009)

Regions	2005										2009									
	Sporting activities	Close to supermarkets	Green space	Public transport	Parking provision	Children's services	Elderly home care	Sporting activities	Close to supermarkets	Green space	Public transport	Parking provision	Children's services	Elderly home care						
Piemonte	34.1	60.3	42.0	189.8	12.5	28.6	1.8	34.1	69.0	42.5	199.3	17.1	37.1	2.3						
Valle d'Aosta	33.9	52.7	23.2	544.0	5.3	100.0	0.1	46.3	58.6	26.2	580.0	8.4	78.4	0.4						
Lombardia	37.7	69.9	27.6	230.1	20.0	54.6	3.2	36.5	68.9	28.6	227.7	24.1	62.5	4.1						
Trentino-Alto Adige	53.1	72.2	71.2	190.7	28.6	75.8	0.6	48.2	71.9	70.3	192.9	34.5	83.8	0.8						
Veneto	39.4	65.8	58.7	122.5	39.8	42.7	5.0	39.6	70.1	62.3	124.4	42.2	70.2	4.8						
Friuli-Venezia Giulia	36.7	72.7	21.8	257.4	11.9	53.0	7.9	37.5	74.6	22.1	258.1	12.0	83.6	7.7						
Liguria	26.6	67.9	35.3	312.5	23.1	75.3	3.1	27.6	70.6	35.4	311.0	22.3	64.3	3.4						
Emilia-Romagna	32.4	71.1	158.5	81.0	24.4	78.0	5.4	36.8	69.3	157.7	83.0	24.0	88.0	8.3						
Toscana	30.4	68.7	152.5	106.0	18.6	78.0	2.1	33.1	64.3	152.1	108.4	20.9	74.6	2.2						
Umbria	31.2	65.9	192.1	162.4	27.4	51.1	4.1	32.3	73.7	187.6	162.8	26.9	63.0	7.6						
Marche	31.4	76.0	185.8	157.2	9.2	45.9	3.3	32.2	67.4	186.1	157.7	15.3	55.7	3.6						
Lazio	33.7	74.3	127.4	124.5	6.5	30.4	3.3	29.4	74.7	121.0	132.3	7.0	30.7	4.0						
Abruzzo	28.9	55.5	714.5	93.5	5.3	26.2	1.8	31.0	63.0	710.0	93.5	21.1	52.1	4.8						
Molise	23.2	52.1	18.3	177.2	1.3	2.9	6.1	22.0	58.7	18.5	177.2	1.2	7.4	2.4						
Campania	22.3	59.3	24.8	227.3	7.3	39.2	1.4	21.1	60.0	25.9	218.0	5.9	50.5	1.9						
Puglia	25.8	70.3	7.8	114.3	7.3	27.5	2.0	23.8	69.6	8.1	122.0	8.2	44.2	2.0						
Basilicata	24.4	55.5	547.9	84.9	2.4	32.8	3.9	27.1	65.2	545.6	87.4	2.3	21.4	5.1						
Calabria	24.5	55.1	19.7	159.6	20.3	7.8	1.6	24.8	56.4	20.8	172.8	19.5	15.6	2.5						
Sicilia	21.5	63.6	71.5	72.2	3.4	33.3	0.8	22.5	68.6	73.3	75.7	6.5	34.6	1.1						
Sardegna	31.1	75.9	86.4	55.7	16.8	17.2	1.1	28.2	78.3	85.9	56.6	16.9	20.4	2.3						
Italia	31.3	67.1	93.5	118.8	14.4	42.8	2.9	31.1	68.5	93.6	122.1	16.2	51.7	3.6						

Source: <http://www3.istat.it/ambiente/contexto/infoterr/assi/asse V.xls>

Table 2 Static and dynamic well-being index (years 2005, 2009)

Region	SWI ₀₅	SWI ₀₉	SWI ₀₉ -SWI ₀₅	DWI _{09/05}
Piemonte	82.1	87.5	5.4	115.6
Valle d'Aosta	63.8	75.5	11.7	128.4
Lombardia	105.0	104.5	-0.5	107.9
Trentino-Alto Adige	106.9	105.1	-1.8	106.7
Veneto	120.9	122.5	1.6	110.0
Friuli-Venezia Giulia	108.6	107.5	-1.1	107.4
Liguria	114.7	105.3	-9.4	99.6
Emilia-Romagna	132.7	134.3	1.7	109.9
Toscana	113.4	107.3	-6.2	102.6
Umbria	136.6	143.6	7.0	114.1
Marche	113.1	115.0	1.9	110.3
Lazio	93.4	87.8	-5.6	102.0
Abruzzo	93.5	137.3	43.8	159.3
Molise	41.6	38.5	-3.1	100.5
Campania	68.4	65.9	-2.5	104.5
Puglia	55.3	55.6	0.3	109.1
Basilicata	89.7	83.7	-6.0	101.3
Calabria	59.5	65.7	6.2	119.8
Sicilia	54.9	60.0	5.1	118.5
Sardegna	70.4	73.6	3.3	113.5

performances of the other regions which has produced a large increase of the national average in 2009.

Overall, the region in which it is possible to record the highest absolute and relative increase of the well-being indicators, over the 5 years, is Abruzzo ($SWI_{09} - SWI_{05} = +43.8$; $DCI_{09/05} = 159.3$); while the largest decrease is observed in Liguria ($SWI_{09} - SWI_{05} = -9.4$; $DCI_{09/05} = 99.6$).

A scatterplot of SWI_{09} against $DWI_{09/05}$ is shown in Fig. 1. Note that the 'static' well-being index may be viewed as the 'speed' of a region and the 'dynamic' well-being index as its 'acceleration.'

In this perspective, we may identify some group of regions by combining different bands of 'speed' and 'acceleration.'

So, Umbria and Emilia-Romagna have high 'speed' and low 'acceleration,' whereas Abruzzo has both high 'speed' and high 'acceleration.'

Vice versa, Molise and Liguria have low 'acceleration' (Liguria goes in reverse since 2005), but Liguria has a very greater 'speed' than Molise.

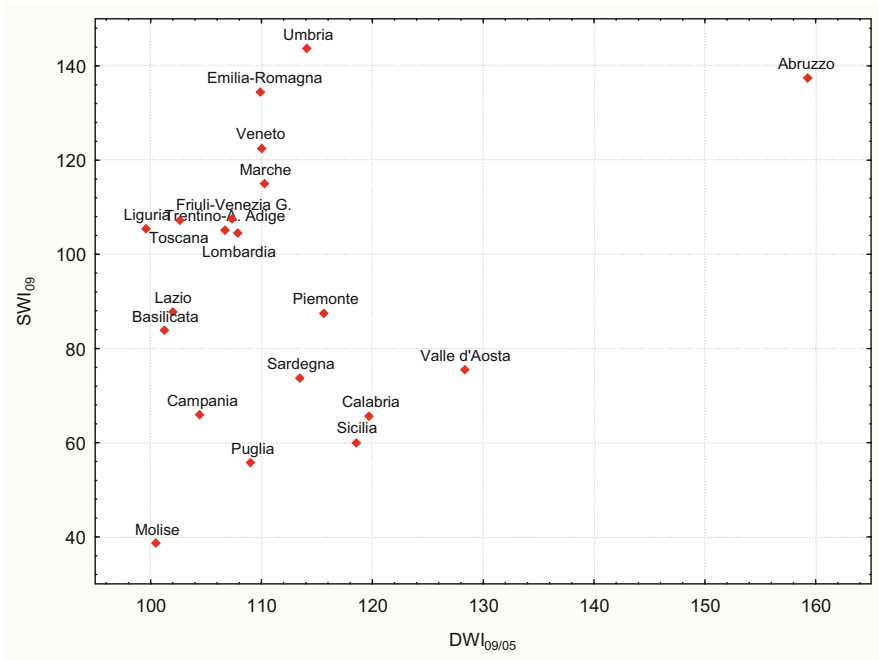


Fig. 1 Comparing static and dynamic well-being index

4 Conclusions

The comparability of the data over time is a central issue in composite indices construction. Normalization methods such as ranking and standardization (z -scores) allow relative comparisons only. Re-scaling in range (0, 1) and indicization by ‘distance to a reference’ measures allow to evaluate absolute changes when the limits or the reference value is independent from the observed data. Another factor that may affect the comparability is the aggregation method, e.g., when the weights are based on Principal Components Analysis or Factor Analysis.

In this paper we propose a method based on the indices number properties for constructing two consistent indices of well-being: a ‘static’ index for spatial comparisons and a ‘dynamic’ index for temporal comparisons.

The method is based on a multiplicative approach and may be applied to different domains without loss of comparability. For example, it is possible to compute the indices for gender and compare the values obtained with other domains.

Acknowledgements Matteo Mazziotta has written Sects. 1 and 3, Adriano Pareto has written Sects. 2 and 4.

References

1. Bandura, R.: A survey of composite indices measuring country performance: 2008 Update. UNDP/ODS Working Paper. Office of Development Studies UNDP, New York (2008)
2. Booyesen, F.: An overview and evaluation of composite indices of development. *Soc. Indic. Res.* **59**, 115–151 (2002)
3. Casadio Tarabusi, E., Guarini, G.: An unbalance adjustment method for development indicators. *Soc. Indic. Res.* **112**, 19–45 (2012). doi:[10.1007/s11205-012-0070-4](https://doi.org/10.1007/s11205-012-0070-4)
4. De Muro, P., Mazziotta, M., Pareto, A.: Composite indices of development and poverty: an application to MDGs. *Soc. Indic. Res.* **104**, 1–18 (2011)
5. Diewert, W.E.: Axiomatic and economic approaches to elementary price indexes. NBER Working Paper Series, 5104. National Bureau of Economic Research, Cambridge (1995)
6. Giovannini, E., Rondinella, T.: Measuring equitable and sustainable well-being in Italy. In: Maggino, F., Nuvolati, G. (eds.) *Quality of Life in Italy: Research and Reflections*, pp. 9–25. Springer, Heidelberg (2012)
7. Michalos, A.C., Smale, B., Labonté, R., Muharjarine, N., Scott, K., Moore, K., Swystun, L., Holden, B., Bernardin, H., Dunning, B., Graham, P., Guhn, M., Gadermann, A.M., Zumbo, B.D., Morgan, A., Brooker, A.-S., Hyman, I.: *The Canadian index of wellbeing. Technical Report 1.0. Canadian Index of Wellbeing and University of Waterloo, Waterloo* (2011)
8. OECD: *Handbook on constructing composite indicators. Methodology and user guide*. OECD Publications, Paris (2008)
9. Tarantola, S.: *European Innovation SCOREBOARD: strategies to measure country progress over time*. JRC Scientific and Technical Reports, EUR 23526 EN. Office for Official Publications of the European Communities, Luxembourg (2008)
10. UNDP: *Human Development Report 2011. Sustainability and equity: a better future for all*. Palgrave Macmillan, New York (2011)
11. Zhou, P., Ang, B.W.: Comparing MCDA aggregation methods in constructing composite indicators using the Shannon-Spearman measure. *Soc. Indic. Res.* **94**, 83–96 (2009)
12. Zhou, P., Ang, B.W., Poh, K.L.: Comparing aggregating methods for constructing the composite environmental index: an objective measure. *Ecol. Econ.* **59**, 305–311 (2006)

Part V

Economic Statistics and Econometrics

A Comparison of Different Procedures for Combining High-Dimensional Multivariate Volatility Forecasts

Alessandra Amendola and Giuseppe Storti

Abstract

The paper investigates the effect of model uncertainty on multivariate volatility prediction. Our aim is twofold. First, by means of a Monte Carlo simulation, we assess the accuracy of different techniques in estimating the combination weights assigned to each candidate model. Second, in order to investigate the economic profitability of forecast combination, we present the results of an application to the optimization of a portfolio of the US stock returns. Our main finding is that, for both real and simulated data, the results are highly sensitive not only to the choice of the model but also to the specific combination procedure being used.

1 Introduction

Due to the variety of different models proposed, model uncertainty is a relevant problem in multivariate volatility prediction. The risk of model misspecification is particularly sizeable in large dimensional problems. In this setting, it is well known that the need for reducing the number of parameters usually requires the formulation of highly restrictive assumptions on the volatility dynamics that, in most cases, are applied without any prior testing (see, e.g., Pesaran et al. [13]). Despite the undoubted relevance of this issue, the analysis of the statistical implications of model uncertainty in multivariate volatility modelling has been largely left unexplored by the statistical and econometric literature. Some recent papers have focused on the evaluation of forecast accuracy of Multivariate GARCH

A. Amendola • G. Storti (✉)

Department of Economics and Statistics (DiSES) & Statlab, University of Salerno, Via Ponte Don Melillo, 84084 Fisciano (SA), Italy

e-mail: alamendola@unisa.it; storti@unisa.it

(MGARCH) models [10–12]. In these papers the main issue of interest has been the derivation of the theoretical properties of some matrix-variate loss functions that can be potentially used for assessing the accuracy of multivariate volatility forecasts. In particular, the attention has been focused on the identification of a class of robust loss functions. In this setting robustness implies that, using a given loss functions for evaluating predictive accuracy, we are able to obtain a ranking of the competing forecasts which is robust to the presence of noise in the volatility proxy used for assessing the forecast accuracy. In their paper Laurent et al. [10] present an empirical comparison of several multivariate volatility forecasts from a set of different MGARCH models using the MCS approach [9] to identify the set of most accurate models. Their results suggest that, independently from market conditions, it is not possible to identify a single model outperforming all the others. This result leaves space for the application of forecast combinations as a tool for improving the forecast accuracy of single, potentially misspecified, volatility models. In a univariate setting, Amendola and Storti [1] have proposed a GMM procedure for the combination of volatility forecasts from different GARCH models. This procedure has been generalized to the combination of multivariate volatility forecasts by Amendola and Storti [2]. The aim of this work is twofold. First, we introduce and compare some alternative combination strategies for multivariate volatility forecasts. All the procedures considered in the paper are not affected by the curse of dimensionality, typically arising in multivariate volatility prediction, and can potentially be applied to the combination of multivariate volatility forecasts for vast dimensional portfolios. In order to assess the accuracy of different combination techniques we present the results of a Monte Carlo simulation study. In particular, our attention will be focused on the analysis of the statistical properties of different estimators of the combination weights. Second, we are interested in assessing the economic profitability of forecast combination in this particular setting. To this purpose, we apply different models and combination strategies to the optimization of a portfolio including a set of stocks traded on the NYSE and compare the performances of the implied optimal portfolios in terms of their volatility.

The paper is structured as follows. In Sect. 2 we will introduce the problem of combining multivariate volatility forecasts from different predictors and discuss some alternative combination strategies. Section 3 will introduce some alternative estimators of the combination weights that can be potentially used in practical applications. The Monte Carlo simulation study will be described in Sect. 4 while the results of the empirical application to stock market data will be presented in Sect. 5. Section 6 will conclude.

2 The Reference Model

The data generating process (DGP) is assumed to be given by

$$\mathbf{r}_t = \mathbf{S}_t \mathbf{z}_t \quad t = 1, \dots, T, T + 1, \dots, T + N \quad (1)$$

where T is the end of the in-sample period, \mathbf{S}_t is any $(k \times k)$ positive definite (p.d.) matrix such that $\mathbf{S}_t \mathbf{S}'_t = \mathbf{H}_t = \text{Var}(\mathbf{r}_t | \mathbf{I}^{t-1})$, $\mathbf{H}_t = C(\mathbf{H}_{1,t}, \dots, \mathbf{H}_{n,t}; \mathbf{w})$ with $\mathbf{H}_{j,t}$ being a symmetric p.d. $(k \times k)$ matrix. In practice $\mathbf{H}_{j,t}$ is a conditional covariance matrix forecast generated by a given “candidate model”. The function $C(\cdot)$ is an appropriately chosen *combination function* and \mathbf{w} is a vector of combination parameters. The weights assigned to each candidate model depend on the values of the elements of \mathbf{w} but do not necessarily coincide with them.

Among all the possible choices of $C(\cdot)$, the most common is the *linear combination function*

$$\mathbf{H}_t = w_1 \mathbf{H}_{1,t} + \dots + w_n \mathbf{H}_{n,t} \quad w_j \geq 0$$

where \mathbf{w} coincides with the vector of combination weights. Alternatively in order to get rid of the positivity constraint on the w_j other combination functions could be selected: the *exponential* and *square root* combination function. The exponential combination is defined as

$$\mathbf{H}_t = \text{Expm} [w_1 \text{Logm}(\mathbf{H}_{1,t}) + \dots + w_n \text{Logm}(\mathbf{H}_{n,t})]$$

where $\text{Expm}(\cdot)$ and $\text{Logm}(\cdot)$ indicate matrix exponential and logarithm, respectively. Differently from the other two functions, the square root combination (for \mathbf{S}_t) is not directly performed on the $H_{j,t}$ but on the $S_{j,t}$

$$\mathbf{S}_t = w_1 \mathbf{S}_{1,t} + \dots + w_n \mathbf{S}_{n,t}$$

with $\mathbf{H}_t = \mathbf{S}_t \mathbf{S}'_t$ and $\mathbf{H}_{j,t} = \mathbf{S}_{j,t} \mathbf{S}'_{j,t}$. In this paper we focus on linear combination functions leaving the investigation of the non-linear combination schemes for future research.

3 Weights Estimators

For the estimation of the combination parameters we consider three different estimation approaches: Composite Quasi ML (CQML), Composite GMM (CGMM) and “pooled” Mincer–Zarnowitz (MZ) regressions. All the estimators considered share the following features: (1) do not imply any assumption on the conditional distribution of returns and (2) can be applied to large dimensional problems. In the CQML method the estimated combination weights are obtained extending to the model in (1) the estimation procedure proposed by Engle et al. [7] in the context of parameter estimation of large dimensional MGARCH models. Formally, the weights w_i are estimated by performing the following optimization:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmax}} \sum_{i \neq j} L(\mathbf{r}^{(ij)} | \mathbf{w}, \mathbf{I}^N),$$

where $\mathbf{r}_t^{(ij)} = (r_{i,t}, r_{j,t})'$, $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_k)'$ and

$$L(\mathbf{r}^{(ij)}|\mathbf{w}, \mathbf{I}^N) = -0.5 \sum_{h=1}^N \log(|\mathbf{H}_{T+h}^{(ij)}|) - 0.5 \sum_{h=1}^N \mathbf{r}_{T+h}^{(ij)} \mathbf{H}_{T+h}^{(ij)} (\mathbf{r}_{T+h}^{(ij)})'$$

is the (bivariate) quasi log-likelihood for the couple of assets (i, j) computed over the prediction period $[T + 1, T + N]$.

In the CGMM estimator the idea is to reduce the problem dimension using the same framework of the CQML. The objective function for estimation is given by the sum of bivariate GMM loss functions, defined as in Amendola and Storti [2], referring to all the feasible bivariate systems that can be extracted from the available dataset. The \hat{w}_i are obtained by performing the following optimization:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i \neq j} m(\mathbf{r}^{(i,j)}; \mathbf{w})' \boldsymbol{\Omega}_N^{(i,j)} m(\mathbf{r}^{(i,j)}; \mathbf{w})$$

where

- $\mathbf{r}_t^{(i,j)} = (r_{i,t}, r_{j,t})$, for $t = T + 1, \dots, T + N$.
- $m(\mathbf{r}^{(i,j)}; \mathbf{w}) = \frac{1}{N} \sum_{t=T+1}^{T+N} \mu(\mathbf{r}_t^{(i,j)}; \mathbf{w})$ and $\mu(\mathbf{r}_t^{(i,j)}; \mathbf{w})$ is a $(p \times 1)$ vector of moment conditions computed from the bivariate system including assets i and j .
- $\boldsymbol{\Omega}_N^{(i,j)}$ is a consistent p.d. estimator of

$$\boldsymbol{\Omega}^{(i,j)} = \lim_{N \rightarrow \infty} NE(m(\mathbf{r}^{(i,j)}; \mathbf{w}^*)m(\mathbf{r}^{(i,j)}; \mathbf{w}^*)')$$

with \mathbf{w}^* being the solution to the moment conditions, i.e. $E(m(\mathbf{r}^{(i,j)}; \mathbf{w}^*)) = \mathbf{0}$.

Finally, as a benchmark for comparison, we consider estimating the combination weights by a “pooled” Mincer–Zarnowitz regression. In this case the \hat{w}_i are given by the OLS estimates of the parameters of the pooled regression model

$$\operatorname{vech}(\tilde{\boldsymbol{\Sigma}}_{T+h}) = w_1 \operatorname{vech}(\tilde{\mathbf{H}}_{1,T+h}) + \dots + w_n \operatorname{vech}(\tilde{\mathbf{H}}_{n,T+h}) + \mathbf{e}_{T+h}$$

for $h = 1, \dots, N$, where depending on the type of combination chosen, $\tilde{\boldsymbol{\Sigma}}_t$ and $\tilde{\mathbf{H}}_{i,t}$ are appropriate transformations of $\mathbf{H}_{i,t}$ and $\boldsymbol{\Sigma}_t$ which is the outer product of returns

$$\boldsymbol{\Sigma}_t = \mathbf{r}_t \mathbf{r}_t'$$

4 Monte Carlo Simulation

In this section we present the results of a simulation study aimed at assessing the accuracy of different estimators of the combination weights assigned to each candidate model. In the simulation we have adopted a linear combination scheme considering, as candidate models, a DCC [5] and a BEKK [6] model with scalar parameters in the dynamic updating equation for conditional correlations and covariances, respectively. Namely, we assume that

$$\mathbf{r}_t = (w_1 \mathbf{H}_{1,t} + w_2 \mathbf{H}_{2,t})^{1/2} \mathbf{z}_t$$

where $\mathbf{z}_t \stackrel{iid}{\sim} MST(5; I_k)$. The cross-sectional dimension of the process has been set equal to $k = 10$. The models for $\mathbf{H}_{1,t}$ and $\mathbf{H}_{2,t}$ are, respectively, given by the following DCC

$$\begin{aligned} \mathbf{H}_{1,t} &= \mathbf{D}_t \mathbf{R}_t \mathbf{D}_t \\ \mathbf{D}_t &= \text{diag}(\mathbf{h}_t) \quad \mathbf{h}_{i,t} = \sqrt{H_{1,ii,t}} \\ H_{1,ii,t} &= a_{0,i} + a_{1,i} r_{i,t-1}^2 + b_{1,i} H_{1,ii,t-1} \\ \mathbf{R}_t &= (\text{diag}(\mathbf{Q}_t))^{-1} \mathbf{Q}_t (\text{diag}(\mathbf{Q}_t))^{-1} \\ \mathbf{Q}_t &= (1 - a - b) \mathbf{R} + a(\boldsymbol{\epsilon}_{t-1} \boldsymbol{\epsilon}'_{t-1}) + b \mathbf{Q}_{t-1} \end{aligned}$$

with $\boldsymbol{\epsilon}_t = \mathbf{D}_t^{-1} \mathbf{r}_t$, and by the scalar BEKK

$$\text{vech}(\mathbf{H}_{2,t}) = (1 - \alpha - \beta) \mathbf{H} + \alpha \mathbf{r}_{t-1} \mathbf{r}'_{t-1} + \beta \mathbf{H}_{2,t-1}$$

where we have set $a = \alpha = 0.03$, $b = \beta = 0.96$ and $\mathbf{H} = \mathbf{d}(\mathbf{H}) \mathbf{R} \mathbf{d}(\mathbf{H})$, with $\mathbf{d}(\mathbf{H}) = \text{diag}(\mathbf{H})^{(1/2)}$. In order to test the robustness of the statistical properties of the weights estimators with respect to the choice of the candidate models, we have constrained the two models to be characterized by equal persistence, which is indeed a critical situation for the implementation of forecasts combination strategies.

We have considered three different time series lengths $N \in \{500, 1000, 2000\}$ and 500 Monte Carlo replicates ($n_{mc} = 500$). As far as the choice of the combination weights w_j is concerned, we consider two different settings. In the first, both candidate models are included into the DGP with weights given by $w_1 = 0.35(0.65)$ and $w_2 = 0.65(0.35)$. In the second setting, only one model is contributing to the DGP which actually means that the weight of the excluded model is set to 0: $w_1 = 0(1)$ and $w_2 = 1(0)$.¹ The simulation results have been summarized by means of box-plots in Figs. 1 and 2. For the sake of brevity, in order to avoid duplication of

¹We remark that, although the DGPs considered in the simulation study do impose a convexity constraint on the combination weights, we do not impose this constraint at the estimation stage.

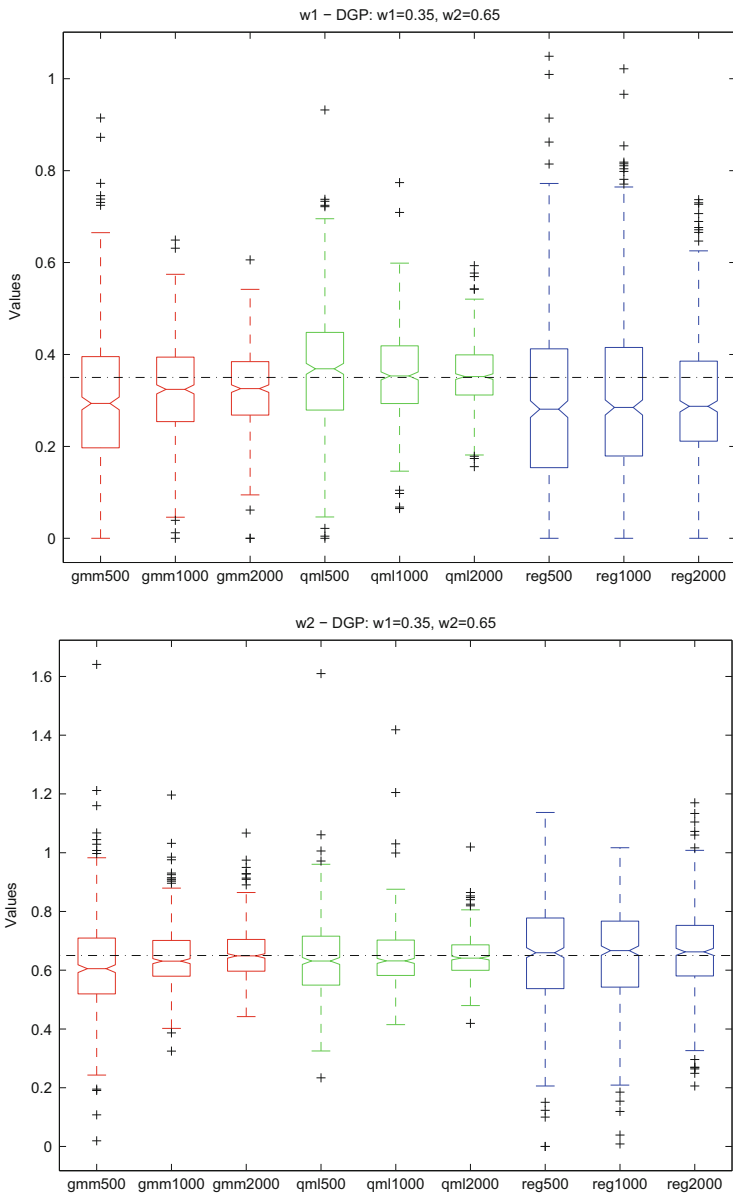


Fig. 1 Simulated distributions of \hat{w}_1 (top) and \hat{w}_2 (bottom). DGP: $w_1 = 0.35, w_2 = 0.65$

information, we only report results for $(w_1 = 0.35, w_2 = 0.65)$ and $(w_1 = 0, w_2 = 1)$. Results for $(w_1 = 0.65, w_2 = 0.35)$ and $(w_1 = 1, w_2 = 0)$ are very similar to those reported and are available upon request. In both the settings considered the CGMM and CQML estimator significantly outperform the estimator

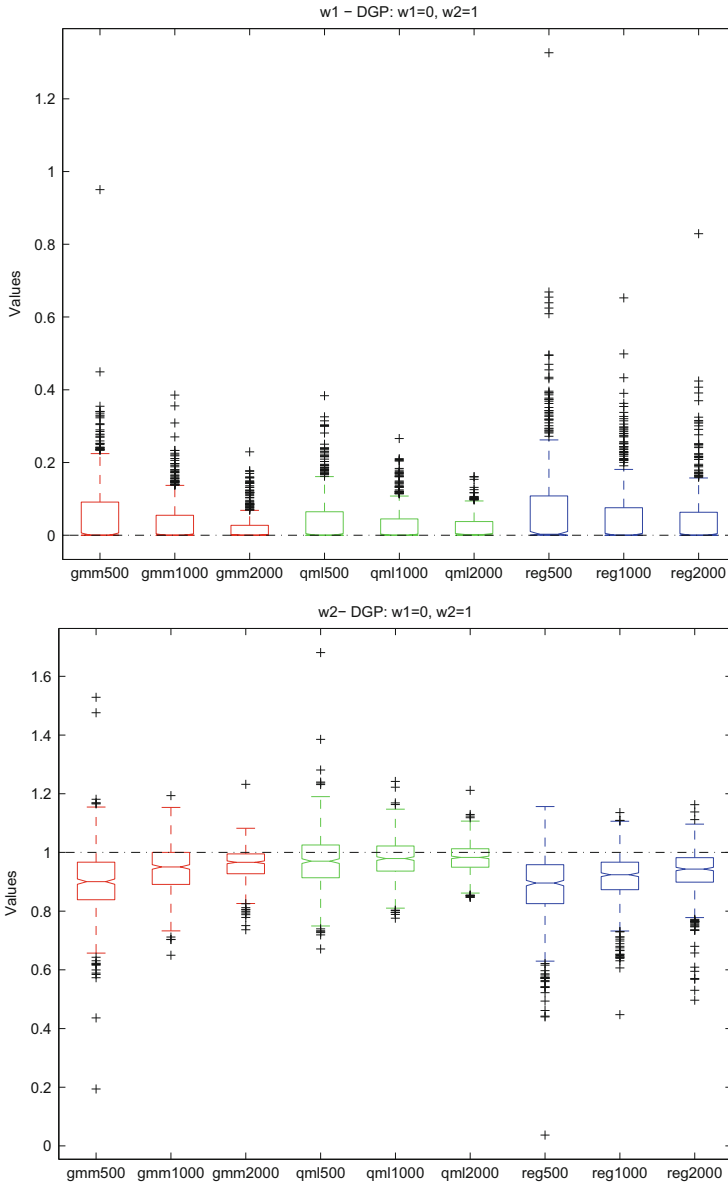


Fig. 2 Simulated distributions of \hat{w}_1 (top) and \hat{w}_2 (bottom). DGP: $w_1 = 0, w_2 = 1$

based on pooled regression with the CQML estimator being slightly more efficient than CGMM. For $(w_1 = 0.35, w_2 = 0.65)$ and $N \leq 1000$, the CGMM estimator is affected by a slight bias component which anyway disappears for $N = 2000$. Differently, the regression based estimator remains biased for any sample size. In

the second setting considered, for $(w_1 = 0, w_2 = 1)$, none of the estimators appears to be significantly biased.

5 An Application to Stock Market Data

In this section we compare the performances of the different combination procedures by means of an application to portfolio optimization. At the same time we want to investigate the economic value of applying forecast combination techniques for determining the asset allocation of a given portfolio. The data we use have been already analyzed by Chiriac and Voev [4].² The dataset contains daily open-to-close returns and realized covariance matrices for six NYSE stocks: American Express Inc. (AXP), Citigroup (C), General Electric (GE), Home Depot Inc. (HD), International Business Machines (IBM) and JPMorgan Chase & Co. (JPM). The sample period starts at January 3, 2000, and ends on March 31, 2011, covering 2828 trading days. As candidate models we consider a scalar DCC, a scalar BEKK, a Constant Conditional Correlation (CCC) model [3],

$$\mathbf{H}_t = \mathbf{D}_t \mathbf{R} \mathbf{D}_t',$$

a RiskMetrics Smoother (ES)

$$\mathbf{H}_t = 0.94\mathbf{H}_{t-1} + (1 - 0.94)\mathbf{r}_t \mathbf{r}_t'$$

and a simple m -terms Moving Covariance (MC) estimator implemented under two different assumptions on the length of the moving window used for volatility prediction, $m = 22$ and 100 days,

$$\mathbf{H}_t = m^{-1} \sum_{j=1}^m (\mathbf{r}_{t-j} \mathbf{r}_{t-j}') \quad m = \{22, 100\}.$$

Our forecasting exercise is based on a rolling-window strategy. We select the first 1500 data point as in-sample period and keep the length of the estimation window fixed over the out-of-sample period. This includes the last 1328 observations that have been used for portfolio optimization and performance evaluation. The estimates of the DCC and BEKK parameters obtained from in-sample data have been reported in Table 1.³

In order to facilitate the comparison between the optimal portfolios yielded by the different techniques considered, the optimization has been performed considering a

²The data can be freely downloaded from the online data archive of the *Journal of Applied Econometrics*. The same data are also used in the paper by Golosnoy et al. [8].

³Due to space constraints, we omit reporting the estimates of the elements of the conditional correlation matrix for the CCC model but this will be made available upon request.

Table 1 Estimates of the parameters of DCC and BEKK models over the in-sample period

	DCC		BEKK
<i>a</i>	0.0046 (0.0003)	α	0.0266 (0.0004)
<i>b</i>	0.9899 (0.0001)	β	0.9713 (0.0004)

Table 2 Estimated combination weights for different models and combination techniques over the in-sample period

Model	REG	REG(rv)	CQML	CGMM
DCC	0.0000	0.0000	0.0000	0.0000
CC	0.0000	0.0000	0.2091	0.3492
ES	0.0000	0.0000	0.4558	0.6218
MCOV(22)	0.0000	0.0000	0.0000	0.0000
MCOV(100)	0.1015	0.1016	0.1638	0.0000
SBEKK	0.5060	0.1961	0.1653	0.1825

target return equal to 0 which amounts to computing the minimum variance portfolio implied by each model. As in the Monte Carlo simulation, we adopt a linear combination strategy while, for estimating the weights, in addition to the CGMM, CQML and pooled regression approach, we also use a pooled regression estimator using realized covariance matrices as dependent variables (REG(rv)). Realized covariances are computed from 5-min returns while subsampling techniques are used to robustify the simple RC estimator with respect to microstructure noise (more details can be found in Chiriac and Voev [4]). The estimated combination weights over the in-sample period have been reported in Table 2. The CQML, CGMM and regression based estimators lead to remarkably different optimal predictors. The CQML and CGMM give similar results with the main difference that the CGMM is excluding the MC(100) estimator. In the regression based approaches the optimal combination only includes the BEKK and MC(100) predictors. Also, for REG and REG(rv), differently from the other two approaches, the sum of weights is much lower than 1 apparently indicating the presence of a substantial negative bias in the candidate predictors.

The optimization performance has been evaluated in terms of the empirical variance of the optimized portfolios over the out-of-sample period (Table 3). As a benchmark we also consider the standard equally weighted (EW) predictor assigning the same weight to each candidate model. What is evident is that the portfolio variance implied by the combined predictors is always lower than that of the candidate models providing evidence in favour of a greater accuracy of the combined predictors as well as of their potential economic profitability for risk management. The minimum variance is obtained when the combination weights are estimated by the CGMM approach. This is probably due to the fact that the CGMM estimator explicitly constraints the covariance matrix of standardized residuals to match, as closely as possible, an identity matrix.

Table 3 Empirical variances of the optimized portfolios

Model	Portfolio variance ^a
DCC	2.33188
CC	2.37658
ES	2.33857
MCOV(22)	2.67185
MCOV(100)	2.10778
SBEEK	2.09339
REG	2.08733
REG(rv)	2.08441
CGMM	2.07337^b
CQML	2.10192
EW	2.08147

^a $\times 10^4$; ^b minimum variance in bold

6 Concluding Remarks

In this paper we have compared different combination techniques for multivariate volatility forecasts by means of a Monte Carlo simulation and an application to portfolio optimization over the US stock market data. The Monte Carlo simulations suggest that the CQML estimator of the combination weights is more efficient than CGMM and pooled regression for all the sample sizes considered. In particular, the CGMM is biased and less efficient than CQML for short sample sizes although its performance improves for $N = 2000$. Differently, in the empirical application the CGMM combination gives the best results in terms of the variance of the optimized portfolios. Our research on the combination of multivariate volatility forecasts is still in progress. Projects for future research include considering applications to large dimensional problems (number of assets ≥ 50), analyzing non-linear combination schemes and a more extensive investigation of the effectiveness of combined volatility predictors in real financial applications.

References

1. Amendola, A., Storti, G.: A GMM procedure for combining volatility forecasts. *Comput. Stat. Data Anal.* **52**(6), 3047–3060 (2008)
2. Amendola, A., Storti, G.: Combination of multivariate volatility forecasts. SFB 649 Discussion Papers, SFB649 DP2009-007, SFB 649, Humboldt University, Berlin (2009)
3. Bollerslev, T.: Modelling the coherence in short-run nominal exchange rates: a multivariate generalized ARCH model. *Rev. Econ. Stat.* **72**(3), 498–505 (1990)
4. Chiriac, R., Voev, V.: Modelling and forecasting multivariate realized volatility. *J. Appl. Econ.* **26**(6), 922–947 (2011)
5. Engle, R.F.: Dynamic conditional correlation: a simple class of multivariate generalized autoregressive conditional heteroskedasticity models. *J. Bus. Econ. Stat.* **20**(3), 339–350 (2002)
6. Engle, R.F., Kroner, K.F.: Modelling the coherence in short-run nominal exchange rates: a multivariate generalized ARCH model. *Econ. Theor.* **11**(1), 122–150 (1995)

7. Engle, R.F., Shephard, N., Sheppard, K.: Fitting vast dimensional time-varying covariance models. Economics Series Working Papers 403, University of Oxford, Department of Economics (2008)
8. Golosnoy, V., Gribisch, B., Liesenfeld, R.: The conditional autoregressive Wishart model for multivariate stock market volatility. *J. Econ.* **167**, 211–223 (2011)
9. Hansen, P.R., Lunde, A., Nason, J.M.: The model confidence set. *Econometrica* **79**, 453–497 (2011)
10. Laurent, S., Rombouts, J.V.K., Violante, F.: On the forecasting accuracy of multivariate GARCH models. *J. Appl. Econ.* **27**(6), 934–955 (2012)
11. Laurent, S., Rombouts, J.V.K., Violante, F.: On loss functions and ranking forecasting performances of multivariate volatility models. *J. Econ.* **173**(1), 1–10 (2013)
12. Patton, A., Sheppard, K.: Evaluating volatility and correlation forecasts. In: Andersen, T.G., Davis, R.A., Kreiss, J.P., Mikosch, T. (eds.) *Handbook of Financial Time Series*, pp. 801–838. Springer-Verlag, Berlin, Heidelberg (2009)
13. Pesaran, M.H., Schleicher, C., Zaffaroni, P.: Model averaging in risk management with an application to futures markets. *J. Empir. Finance* **16**(2), 280–305 (2009)

Which Seasonality in Italian Daily Electricity Prices? A Study with State Space Models

Paolo Chirico

Abstract

The paper presents a study of seasonality in Italian daily electricity prices. In particular, it compares the ARIMA approach with the structural state space approach in the case of seasonal data. Unlike ARIMA modeling, the structural approach has enabled us to detect, in the prices under consideration, the presence of stochastic daily effects whose intensity is slowly decreasing over time. This dynamic of seasonality is the consequence of a more balanced consumption of electricity over the week. Some causes of this behavior will be discussed in the final considerations. Moreover, it will be proved that state space modeling allows the type of seasonality, stochastic or deterministic, to be tested more efficiently than when unit root tests are used.

1 Introduction

In the past 20 years, competitive wholesale markets of electricity have started in the OECD countries in the international context of the deregulation of energy markets. At the same time, an increasing number of studies on electricity prices have been published. Most of these studies have sought to identify good prediction models, and, for this reason, ARIMA modeling has been the most common methodology. Nevertheless, electricity prices present periodic patterns, seasonality in time series terminology, for which ARIMA modeling does not always seem to be the best approach.

P. Chirico (✉)

Department of Economics and Statistics, University of Turin, Turin, Italy
e-mail: paolo.chirico@unito.it

The treatment of seasonality in the ARIMA framework is conceptually similar to the treatment of trends: like these, seasonality entails the non-stationarity of the process, and its non-stationary effect has to be removed before modeling the process. More specifically, if the seasonal effects are constant at corresponding times (e.g., every Sunday, every Monday, etc.), the seasonality can be represented by a periodic linear function $s(t)$ (*deterministic seasonality*). In this case, the correct treatment consists in subtracting the seasonality, and then in modeling the non-seasonal prices using an ARIMA model¹:

$$\phi(B)\Delta[p_t - s(t)] = \theta(B)\varepsilon_t. \quad (1)$$

On the other hand, if the seasonal effects are characterized by stochastic variability (*stochastic seasonality*), the correct treatment consists in applying the seasonal difference to the prices $\Delta_s p_t = p_t - p_{t-s}$, and then modeling the differences using an ARIMA model:

$$\phi(B)\Delta\Delta_s p_t = \theta(B)\varepsilon_t. \quad (2)$$

The two treatments are not interchangeable. In fact, in the case of deterministic seasonality, the seasonal difference is not efficient because it introduces seasonal unit roots into the moving average part $\theta(B)$ of the ARIMA model; in the case of stochastic seasonality, the first treatment does not assure stationarity in the second moment of the data [4]. Hence, the correct application of ARIMA models to seasonal data requires first the identification of the type, stochastic or deterministic, of the seasonality present in the prices.

In many cases, statistical tests indicate the presence of deterministic seasonality, at least in the short run. For this reason, as well as for easiness reasons, many scholars [1, 7, 10] have opted for representing seasonality by means of periodic functions. This approach makes it possible to measure the seasonal effects, but it is based on the strong assumption that seasonal effects remain constant over time. On the other hand, the seasonal difference approach does not satisfy the need to understand and model the real dynamic of seasonality in electricity prices. Therefore, other scholars [8] turned to periodic ARIMA models, but this modeling requires numerous parameters when seasonality presents numerous periods (e.g., daily pattern). On the basis of these considerations, structural state space models could be a solution for representing seasonality in a flexible way, but using few parameters. The paper illustrates some structural space state models that yielded interesting findings about seasonality in Italian daily electricity prices. More specifically, the paper is organized as follows. The next section illustrates some items about deterministic and stochastic seasonality, and the most common test for checking seasonality is presented. In Sect. 3, an analysis of the Italian daily

¹Formally, model 1 is called the Reg-ARIMA model by some authors, ARMAX by others.

electricity prices is discussed, comparing the ARIMA approach with the structural (space state) approach. Final considerations are made in the last section.

2 Deterministic and Stochastic Seasonality

Seasonality can be viewed as a periodic component s_t of a seasonal process y_t that makes the process non-stationary:

$$y_t = y_t^{\text{ns}} + s_t. \quad (3)$$

The remaining part $y_t^{\text{ns}} = y_t - s_t$ is the non-seasonal process and is generally assumed stochastic, but s_t can be either deterministic or stochastic.

Deterministic seasonality can be represented by periodic functions of time (having s periods) like the following ones:

$$s_t = \sum_{j=1}^s \gamma_j d_{j,t} \text{ with } \sum_{j=1}^s \gamma_j = 0 \quad (4)$$

$$\text{or } s_t = \sum_{j=1}^{\lfloor s/2 \rfloor} A_j \cos(\omega_j t - \phi_j). \quad (5)$$

In Eq. (4), the parameter γ_j represents the seasonal effect in the j -th period ($d_{j,t}$ is a dummy variable indicating the period). In Eq. (5), seasonality is viewed as the sum of $\lfloor s/2 \rfloor^2$ harmonic functions each of them having angular frequency $\omega_j = j2\pi/s$; $j = 1, 2, \dots, \lfloor s/2 \rfloor$. Deterministic seasonality satisfies the following relation:

$$S(B)s_t = 0 \quad (6)$$

where $S(B) = 1 + B + B^2 + \dots + B^{s-1}$ is the seasonal summation operator based on the backward operator B .³ In the case of stochastic seasonality, the relation (6) becomes:

$$S(B)s_t = w_t \quad (7)$$

where w_t is a zero-mean stochastic process (stationary or integrated). Now seasonality can be viewed as the sum of $\lfloor s/2 \rfloor$ stochastic harmonic paths $h_{j,t}$:

$$\gamma_j(B)h_{j,t} = w_{j,t} \quad (8)$$

² $\lfloor s/2 \rfloor = s/2$ for s even, and $\lfloor s/2 \rfloor = (s-1)/2$ for s odd.

³ $S(B)s_t = s_t + s_{t-1} + \dots + s_{t-s+1}$.

where

$$\gamma_j(B) = (1 - e^{i\omega_j}B)(1 - e^{-i\omega_j}B) \text{ if } 0 < \omega_j < \pi \quad (9)$$

$$\gamma_j(B) = (1 + B) \text{ if } \omega_j = \pi \quad (10)$$

and $w_{j,t}$ is a *zero-mean* stochastic process (stationary or integrated).

Since each seasonal operator $\gamma_j(B)$ is a polynomial with unit roots, each stochastic harmonic path implies the presence of one or two (complex and conjugate) unit roots in the process (more exactly, in the autoregressive representation of the process) and vice versa. Finally, since:

$$\Delta_s = \Delta S(B) = \Delta \prod_{j=1}^{\lfloor s/2 \rfloor} \gamma_j(B) \quad (11)$$

the application of the filter Δ_s to a seasonal process y_t makes the process stationary, removing a stochastic trend (eventually present in the non-seasonal data) and $\lfloor s/2 \rfloor$ stochastic harmonic paths present in seasonality.

2.1 HEGY Test

A very common methodology used to test for non-stationarity due to seasonality is the procedure developed by Hylleberg et al. [6], and known as the HEGY test. This test was originally devised for quarterly seasonality, but it has also been extended for weekly seasonality in daily data by Rubia [11].

Under the null hypothesis, the HEGY test assumes that the relevant variable is *seasonally integrated*. This means, in the case of daily electricity prices (p_t), that the weekly difference $\Delta_7 p_t$ is assumed to be a stationary process.

Since:

$$\Delta_7 = (1 - B) \prod_{j=1}^3 (1 - e^{i\omega_j}B)(1 - e^{-i\omega_j}B) \quad (12)$$

($\omega_j = 2\pi/7, 4\pi/7, 6\pi/7$), the null hypothesis of the HEGY test entails the presence in the process of seven unit roots: one at zero frequency (corresponding to a stochastic trend) and three pairs of complex unit roots corresponding to three stochastic harmonic paths with frequencies $2\pi/7, 4\pi/7, 6\pi/7$.

The test consists in checking the presence in the process of the unit roots; in this sense it can be viewed as an extension of the Dickey–Fuller tests [2]. Like these

tests, the HEGY test is based on an auxiliary regression⁴:

$$\Delta_7 p_t = \alpha + \sum_{s=2}^7 \gamma_s d_{s,t} + \sum_{r=1}^7 \alpha_r z_{r,t-1} + \sum_{j=1}^p \phi_j \Delta_7 p_{t-j} + \varepsilon_t \quad (13)$$

where $d_{s,t}$ is a zero/one dummy variable corresponding to the s -th day of the week, and each regressor $z_{r,t}$ is obtained by filtering the process p_t so that:

- it will be orthogonal to the other regressors;
- it will include only one root of the seven roots included in p_t .

For example, $z_{1,t}$ includes only the unit root having zero frequency (stochastic trend), but not the seasonal roots; $z_{2,t}$ and $z_{3,t}$ include only the seasonal roots having frequency $2\pi/7$ and so on (see [11] for more details).

The number p of lags of the dependent variable in the auxiliary regression (augmentation) has to be chosen to avoid serial correlation in the error term ε_t .

If $\Delta_7 p_t$ is a stationary process, all roots have been removed, and the coefficients α_s are not significant. As in the augmented unit root test of Dickey and Fuller (ADF), the null hypothesis $\alpha_1 = 0$ is accepted against the alternative hypothesis $\alpha_1 < 0$ on the basis of a non-standard t -statistic. In regard to the seasonal roots, the test should be performed on each couple of roots having the same frequency. Indeed, only the hypothesis $\alpha_{2j} = \alpha_{2j+1} = 0$ ($k = 1, 2, 3$) means the absence in $\Delta_7 p_t$ (i.e., the presence in p_t) of a harmonic path with frequency $2\pi j/7$. This assumption can be tested by a joint F -test; the distribution of each statistic F_j is not standard, but the critical values are reported in [9]. In conclusion: if some hypothesis $\alpha_{2j} = \alpha_{2j+1} = 0$ is not rejected, the seasonality should be stochastic; if all the hypotheses $\alpha_{2j} = \alpha_{2j+1} = 0$ are rejected and some coefficient γ_s is significant, the seasonality should be deterministic.

3 Analysis of the Italian Daily Electricity Prices

The HEGY test was performed on the 2008–2011 Italian daily PUN⁵ (more specifically the log-PUN). As reported in Table 1, none of the null hypotheses (H_0) was significant at 1% level. Nevertheless, the absence of a stochastic trend was not confirmed by the ADF test on the same data. This might mean that the prices process is nearly a stochastic trend, but also that the process is not homogeneous over the whole period. Indeed, after performing the HEGY test on the sub-periods 2008–2009 and 2010–2011, it can be noted that the statistic t concerning the

⁴This is a standard version of the HEGY test for daily data, but it can be extended to include trends. Nevertheless, in this case, there is no reason for doing so.

⁵The PUN is the National Single Price in the Italian electricity market (IPEX). The PUN series are downloadable from the web site of the Energy Markets Manager: <http://www.mercatoelettrico.org>.

presence of a stochastic trend gives different signals: the 2008–2009 daily prices seem to include a stochastic trend, whereas the 2010–2011 daily prices do not. Such deductions were confirmed by performing the ADF test on the data (Table 1). The absence of mean-reversion in the first period is a particular case and should be related to the high variation of the oil prices in the same period. On the other hand, seasonality remains nonstochastic in both periods (absence of seasonal roots). According to these findings, the 2008–2009 daily log-PUN was represented by a Reg-ARIMA model, but the 2010–2011 log-PUN by a Reg-ARMA model: more specifically, a Reg-IMA(1,2) for the first period and a Reg-AR(7) for the second one. In both cases the regression was the following:

$$p_t = \gamma_{\text{Mon}}d_{\text{Mon},t} + \dots + \gamma_{\text{Sat}}d_{\text{Sat},t} + p_t^{\text{ns}}. \tag{14}$$

The models parameters and their significance are reported in Table 2.

Table 1 HEGY and ADF tests

H_0	Stat.	2008–2011 (sign.)	2008–2009 (sig.)	2010–2011 (sign.)
$\alpha_1 = 0$	t	−3739 ***	−1572	−3742 ***
$\alpha_2 = \alpha_3 = 0$	F_1	163,127 ***	65,480 ***	84,874 ***
$\alpha_4 = \alpha_5 = 0$	F_2	175,639 ***	66,215 ***	90,891 ***
$\alpha_6 = \alpha_7 = 0$	F_3	232,567 ***	103,114 ***	93,018 ***
ADF test	τ	−2288	−1167	−3371**

p-value < 0.05; *p-value < 0.01

Table 2 Models parameters

Param.	Model 2008–2009 value/sign.	Model 2010–2011 value/sign.
Const	−0.001	4.136 ***
AR1	−	0.351 ***
AR2	−	0.116 ***
AR3	−	0.084 **
AR4	−	0.134 ***
AR5	−	0.088 **
AR6	−	0.038
AR7	−	0.067 *
MA1	−0.498 ***	−
MA2	−0.237 ***	−
Mon	0.148 ***	0.076 ***
Tue	0.173 ***	0.093 ***
Wed	0.190 ***	0.091 ***
Thu	0.169 ***	0.097 ***
Fri	0.151 ***	0.080 ***
Sat	0.103 ***	0.072 ***

*p-value < 0.10; **p-value < 0.05;

***p-value < 0.01

Since the analyzed data are log-prices, each daily coefficient (lower part of the table) indicates the average per-cent difference between the corresponding daily price and the Sunday price, which is obviously the lowest price. Indeed, the consumption of electricity is generally lowest on Sundays. To be noted is that the daily effects are lower in the second period. This result may mean that there was a structural break in the seasonality as a consequence of a structural break in the daily demands or in the daily supplies of electricity. On the other hand, seasonality may have had fluctuations of slowly decreasing intensity in the period 2008–2011 as a consequence of slow changes in the daily demands and/or daily supplies of electricity.

In order to gain better understanding of the dynamics of seasonality in the electricity prices, we analyzed the prices by means of state space models.

3.1 State Space Analysis of Electricity Prices

First, the following model was performed on the 2008–2011 daily log-PUN:

$$p_t = m_t + s_t + \varepsilon_t \quad (15)$$

$$m_{t+1} = m_t + b_t + \varepsilon_{1,t} \quad (16)$$

$$b_{t+1} = b_t + \varepsilon_{2,t} \quad (17)$$

$$s_{t+1} = -s_t - s_{t-1} - \dots - s_{t-5} + \varepsilon_{3,t} \quad (18)$$

where m_t is the non-seasonal level of the log-PUN p_t ; b_t is the slope and s_t is the seasonality (daily effect). The disturbance factors ε_t , $\varepsilon_{1,t}$, $\varepsilon_{2,t}$ and $\varepsilon_{3,t}$ are *white noises* with variances σ^2 , σ_1^2 , σ_2^2 and σ_3^2 .

This model, also known as the *local linear trend model with seasonal effect* [3], is a common starting state space model for seasonal data. Equation (18) is a particular case of assumption (7) (w_t is assumed to be a white noise) and entails stochastic seasonality. This assumption permits seasonality to change in the period 2008–2010 according to the findings in Table 1. The estimation results of this model are reported in the second column of Table 3. To be noted is that the estimate of

Table 3 Three state space models for the log-prices

Parameters	Local trend with seasonal effect	Local level with seasonal effect	Local level with decreasing seas. eff.
σ	0.0856 ***	0.0858 ***	0.0858 ***
σ_1	0.0401 ***	0.0397 ***	0.0397 ***
σ_2	0.0000	No	No
σ_3	0.0034 ***	0.0034 ***	0.0029 ***
α	No	No	0.9923
AIC	-2314.61	-2346.51	-2365.15
SBC	-2293.11	-2330.37	-2343.64

***p-value < 0.01

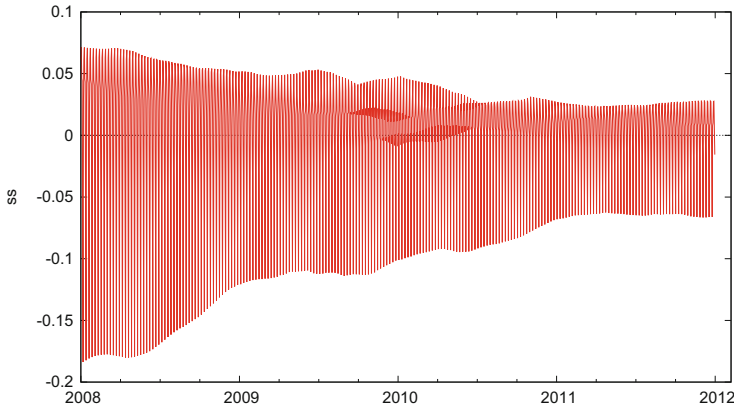


Fig. 1 Seasonality in the period 2008–2011

the standard deviation of ε_2 is zero, which means the slope of the trend b_t can be assumed to be nonstochastic; moreover, the estimate of b_t converges to zero. For these reasons, a seasonal model without slope (without b_t in Eq. (16) and without Eq. (17)), also known as the *local level with seasonal effect*, shows better indices of fit (third column).

The consideration of the diagram of the smoothed seasonality (Fig. 1) shows that the daily effects tend to decrease in the period 2008–2011 (in this case, the daily effects should be viewed as the percentage deviations, positive on working days and negative at weekends, from the trend of prices). According to this evidence, the standard local level model was modified by the following seasonal state equation:

$$s_{t+1} = -\alpha(s_t + s_{t-1} + \dots + s_{t-5}) + \varepsilon_{3,t} \quad (19)$$

where $0 < \alpha < 1$ so that the daily effects can tend to decrease. Performing the new model on the 2008–2011 log-PUN, the value of alpha resulted equal to 0.9923 (Table 3, fourth column); the standard deviation of the disturbance on seasonality was equal to 0.0029 (less than 0.3 %). The values of the Akaike (AIC) and Schwarz (SBC) indices are less than in previous models, denoting an improvement in fit. These findings prove that the daily effects were very slowly decreasing in the period 2008–2011; indeed, so slowly decreasing and so little varying that they could be viewed as constant in a short period. For this reason, the HEGY test, which is not a particularly powerful test, detected deterministic seasonality (Table 1).

4 Final Considerations

The analysis described in the previous sections has shown that the daily effects (i.e., seasonality) on daily wholesale electricity prices exhibited slowly decreasing intensity in the period 2008–2011 in Italy. We reiterate that the daily effects can

be viewed as deviations from the trend of the prices due to the days of the week. A reduction of the daily effects means a reduction of the differences among the daily prices. Some causes regarding the demand and the supply of electricity can be highlighted. In regard to the demand, a more balanced consumption of electricity over the week has been noted in recent years. One reason is certainly that more and more families have subscribed contracts of domestic electricity provision which make electricity consumption cheaper in the evenings and at weekends. Moreover, the difficulties of the Italian economy in recent years have caused a reduction in electricity consumption on working days.

In regard to the supply, the entry into the market of several small electricity producers has made the supply of electricity more flexible.

Regarding the methodology, structural state space models seem to be a more powerful tool than the HEGY test for detecting the type of seasonality. From the state sequence of the seasonal components, it is possible to gain a first view on the kind of seasonality affecting the data. Moreover, the significance test on the standard deviation of the disturbance in the seasonal component makes it possible to check whether or not seasonality is stochastic. More specifically, if the standard deviation is not significant, the seasonality should be assumed to be deterministic; otherwise it should be assumed stochastic. Although these models are not usually employed for electricity prices, they have interesting features for the analysis and prediction of electricity prices. As known, state space modeling can include ARIMA modeling, but it allows easier modeling of periodic components compared with the latter. Moreover, structural state space models can represent electricity prices according to an economic or behavioral theory.

This study has not dealt with volatility clustering, a well-known feature/problem of electricity prices. As known, the GARCH models (in all versions) are typically used to model volatility clustering. Although such modeling is generally associated with ARIMA modeling, conditional heteroscedasticity can be considered in structural framework as well [5].

References

1. Bhanot, K.: Behavior of power prices: implications for the valuation and hedging of financial contracts. *J. Risk* **2**, 43–62 (2000)
2. Dickey, D.A., Fuller, W.A.: Distribution of the estimators for autoregressive time series with a unit root. *J. Am. Stat. Assoc.* **74**, 427–431 (1979)
3. Durbin, J., Koopman, S.J.: *Time Series Analysis by State Space Methods*. Oxford University Press, Oxford (2001)
4. Hamilton, J.D.: *Time Series Analysis*. Princeton University Press, Princeton (1994)
5. Harvey, A., Ruiz, E., Sentana, E.: Unobserved component time series models with ARCH disturbances. *J. Econ.* **52**, 129–157 (1992)
6. Hylleberg, S., Engle, R.F., Granger, C.W.J., Yoo, B.S.: Seasonal integration and cointegration. *J. Econ.* **44**, 215–238 (1990)
7. Knittel, C.R., Roberts, M.R.: An empirical examination of restructured electricity prices. *Energy Econ.* **27**, 791–817 (2005)

8. Koopman, S.J., Ooms, M., Carnero, M.A.: Periodic seasonal Reg-ARFIMA-GARCH models for daily electricity spot prices. *J. Am. Stat. Assoc.* **102**(477), 16–27 (2007)
9. Leon, A., Rubia, A.: Testing for weekly seasonal unit roots in the Spanish power pool. In: Bunn, D.W. (ed.) *Modelling Prices in Competitive Electricity Markets*, pp. 131–145. Wiley, London (2004)
10. Lucia, J., Schwartz, E.: Electricity prices and power derivatives: evidence from the Nordic Power Exchange. *Rev. Deriv. Res.* **5**(1), 5–50 (2002)
11. Rubia, A.: Testing for weekly seasonal unit roots in daily electricity demand: evidence from deregulated markets. Instituto Valenciano de Investigaciones Economicas, WP-2001-21. <http://www.ivie.es/downloads/docs/wpasec/wpasec-2001-21.pdf> (2001)

From the Standard of Living as a Latent Variable to the Estimation of Equivalence Scales and Other Indices

Gustavo De Santis and Mauro Maltagliati

Abstract

A recent approach to the estimation of equivalence scales (S) suggests that a three-step procedure be followed: one must first form clusters of households with the same apparent standard of living (a latent dimension, to be inferred from selected, “well-behaved” indicators), then estimate within-cluster equivalence scales and finally, with a weighted average, obtain the general equivalence scale. This paper further elaborates on these ideas and illustrates how the same logic can also lead to the estimation of inflation and *PPP* (purchasing power parity). Thanks to its flexibility, the method can be applied not only to “standard” databases (expenditure surveys) but also to income surveys (e.g. the Bank of Italy *SHIW*—Survey on Household Income and Wealth), and to any other database including an indicator of resources (e.g. income or total expenditure) and a few “well-behaved” indicators of economic well-being. Empirical results for Italy (2004–2010) are presented and discussed.

1 Equivalence Scales: Economists vs. Statisticians

Ever since the times of Engel [4], (micro)economists have been interested in equivalence scales, index numbers summarising the extra needs of households that somehow differ from the reference one. Unfortunately, after innumerable papers and debates (see, e.g., [8, 11]), “almost every aspect of equivalence scale specification remains controversial” [9], and the gloomy conclusion now seems to be that “in

G. De Santis • M. Maltagliati (✉)

Dip. di Statistica, Informatica, Applicazioni “G. Parenti”, Viale Morgagni 59,
50134 Florence, Italy

e-mail: gustavo.desantis@unifi.it; mmaltag@disia.unifi.it

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,
Studies in Theoretical and Applied Statistics,

DOI 10.1007/978-3-319-27274-0_25

general, there is no accepted method for determining equivalence scales, and no equivalence scale [can be] recommended" [10].

Theoretically elusive as they may be, however, equivalence scales are in practice indispensable (e.g. for taxes and subsidies, or for the analysis of household poverty), so that governments and official entities have resorted to "empirical" and very simple scales that are either scarcely defensible (e.g. in Italy, where Istat and the government use scales that are based on Engel's food share) or totally arbitrary, but politically "strong" (e.g. the OECD scales—which, incidentally, have repeatedly changed over time: the currently prevailing version is the square root one, by which the equivalence scale is assumed to evolve with the square root of the household dimension; [10]).

It is not easy to say why equivalence scales derived from apparently sound microeconomic assumptions (starting, typically, from the so-called complete demand system) do not work: they prove inconsistent, their parameters change abruptly and illogically. This can occur from one year to the next or in the same year and on the same database by a simple reorganisation of expenditure shares, which are the basic elements of this procedure. Some argue that equivalence scales *cannot* work, because of their implicit compensation nature: most changes in the household structure derive from free choices (e.g. to live with a partner, or to have a child), and these, by definition, cannot lower one's standard of living. Nonetheless, an equivalence scale will illogically indicate that larger households need extra resources to be restored to their original level of well-being [12].

But this explanation is not totally convincing, for two main reasons. The first is that inconsistencies do not depend on what one tries to measure, and emerge both with "choices" (e.g. an extra child) and with other types of household change (e.g. getting older, or falling ill). The second reason relates to the very notion of "standard of living": economists assume that consumers make all their decisions simultaneously, so that one of them (e.g. a large family) cannot explain the others (e.g. consumption behaviour). We however believe that some of these decisions come first, or are inherited from the past (e.g. fertility and children), and that, together with other dimensions (resources, but also preferences, culture, climate, etc.), they affect the way households spend their money. This link between expenditure behaviour, resources and household structure is precisely the reason why equivalence scales can be estimated.

The rest of the paper is organised as follows: first the method is introduced, with its rationale (Sect. 2) and its formulae (Sect. 3). This method is then applied to two different datasets in Italy (income and consumption survey, 2004–2010), an approach which is a novelty in itself. In addition, inflation and regional purchasing power parities (*PPP*) are estimated together with equivalence scales, another novel approach (Sect. 4). The two databases yield results that are comparable between them, and consistent with the (scarce) external evidence available (Sect. 5). We conclude that, with this approach, many more databases can be used for the estimation of equivalence scales (and other indices) than was thus far deemed possible (Sect. 6).

2 The Standard of Living as a Latent Variable

Let us define the standard of living of a household as the ratio between its resources (e.g. income, or total expenditure) and its needs (to be measured, in relative terms, by means of an equivalence scale). The standard of living is not directly observable: it is a latent variable that we assume to be associated with (and therefore revealed by) a few, properly selected empirical indicators. A priori, one could draw a list of “classical” indicators of affluence, ranging from the way people dress to the clubs they are members of, from where they do their shopping to the neighbourhoods they live in and so on. In this paper, however, a different path is followed: well-being indicators are not pre-determined, but they are chosen *ex post*, on the basis of actual household expenditure patterns. Our indicators are *all* the variables that satisfy two requirements: firstly, they must vary monotonically with economic resources (income Y), for any given household size N and for all the years of observation, if there is more than one. Secondly, they must not bear any trivial relation with household dimension N . Take total outlays X , for instance: they increase with Y for any given household dimension N , but they are also positively affected by N , given Y , because more people have more needs and increase household expenditure, without any clear connection with per capita expenditure or individual “utility”. Therefore total expenditure X does not qualify as a valid, or “well-behaved” indicator for our purposes.

Of course, indicators that depend solely, or at least primarily, on the standard of living would be preferable, but, in practice, one has to keep all the well-behaved indicators that can be found, hoping that the other variables which influence them (preferences, culture, climate, etc.) will average out and will not significantly affect the results.

The food share, for instance, is a well-behaved indicator of economic well-being (it decreases as income increases, for all household dimensions, in all years), and this is why the Engel method works well, in practice. But the use of a single indicator is dangerous. Apart from possible (indeed, frequent) measurement errors and biases, including those that derive from observing households for only a very short period (as it normally happens with expenditure surveys), almost all the available indicators, including food, are also affected by something else than just economic affluence, which we may call “style”: they depend, among other things, on gender (women may buy lipsticks, and men sportswear), age (toys vs. canes), religious beliefs (pork and alcohol are banned in some cases), environment and season (bathing costumes vs. scarves), health conditions (diabetics do not buy sugar) and so on. The food share, too, is not immune from these risks: for instance, it tends to increase when there are very young and very old members in the households, since in both cases most time is spent at home, and other consumption items matter proportionally less. With several, relatively independent well-behaved indicators of economic well-being this risk is reduced, even though it can never be totally ruled out.

Once a sufficient number of well-behaved indicators have been collected, such indicators are used to form clusters of households. Within each cluster k , by construction, all households are assumed to share a reasonably similar standard of living: we calculate the average income of structurally homogeneous households (for instance, with the same dimension: N , or $N + 1$, $N + 2$, etc.), and we consider this average income $Y_{n,k}$ as representative of the “typical” N -member household of cluster k . Therefore, the ratios $Y_{n+m,k}/Y_{n,k}$ define the equivalence scales for cluster k and the average of these cluster-specific scales yields the general equivalence scale.

An extension of this idea is one of the novel approaches described in this paper. The same procedure can in fact be followed when households, within clusters, are stratified on the basis of different criteria: for instance, by region of residence (which leads to the construction of *PPP*—purchasing power parity—indices) or by year of observation (which leads to an estimate of inflation). Note that neither income Y , nor the stratification variables (dimension, region or year of observation) are used for clustering, which in turn means that they are not used for evaluating the standard of living of households.

3 Cluster-Specific and General Indices

Let us assume households H to be characterised by the following set of variables: \mathbf{E} (a vector of well-behaved economic indicators), Y (income), N (number of members), T (year of observation) and R (region of residence)

$$H = (\mathbf{E}_h, Y_h, N_h, T_h, R_h) \quad (1)$$

Households are clustered on the basis of the \mathbf{E} (observable, well-behaved) variables, while the other variables (Y , N , T and R) will be used at a later stage. Since all \mathbf{E} , by definition, are linked to economic well-being (a latent, unobservable variable), we assume that all the households that fall in the same cluster share a similar level of economic affluence, which is an essential requirement for comparison and for the construction of our indices.

If this assumption holds, the average income Y of all households h within cluster k , with given characteristics (dimension n , in year t , in region r)

$$Y_{k,n,t,r} = \frac{\sum_h c_{h,k,n,t,r} Y_{h,k,n,t,r}}{\sum_h c_{h,k,n,t,r}} \quad (2)$$

gives an estimate of the “true” level of income of these households, where the coefficient c represents the weight of each specific observation, if this is provided in the database (as it normally happens with not representative, e.g. stratified, samples).

Within cluster k , by construction, the (unknown) level of economic well-being is (almost) the same for all households, and this allows us to estimate the cluster-specific indices we are interested in. For instance, for equivalence scales we get

$$S_{k,n,t,r} = \frac{Y_{k,n,t,r}}{Y_{k,n=1,t,r}} \tag{3}$$

(where one-person households, $n = 1$, are used as a standard of reference) and, from these, we can compute the general equivalence scale, as a weighted mean of cluster-specific estimates

$$S_n = \sum_k \sum_t \sum_r S_{k,n,t,r} w_{k,n,t,r} \tag{4}$$

where the weights $w_{k,n,t,r} \left(= \frac{\sqrt{H_{k,n,t,r} \cdot H_{k,n=1,t,r}}}{\sum_k \sum_t \sum_r \sqrt{H_{k,n,t,r} \cdot H_{k,n=1,t,r}}} \right)$ take into account the number of households H both in the numerator and in the denominator of each of these ratios [3].

In addition, we can also compute the variance of the estimated indices, e.g. of equivalence scales. Dropping a few of our classification letters for the sake of simplicity, within clusters, the variance of S is [7]

$$Var (S_{k,n}) = \frac{1}{(Y_{k,n=1})^2} \{Var (Y_{k,n}) + S_{k,n}^2 [Var (Y_{k,n=1})]\} \tag{5}$$

where $Var(Y_{k,n})$ stands for the estimated variance of the average income, for any given combination of cluster k and household dimension n .

The estimated variance of S_n is [7]

$$Var (S_n) = \sum_k w_{k,n}^2 \cdot Var (S_{k,n}) \tag{6}$$

If instead it is PPP that we want to estimate, our basic elements (in 3) are to be constructed in a slightly different manner

$$PPP_{k,n,t,r} = S_{k,n,t,r} = \frac{Y_{k,n,t,r}}{Y_{k,n,t,r=1}} \tag{7}$$

where an arbitrary region (e.g. the North-West) is taken as a standard of reference. Formula (4) then becomes

$$PPP_r = S_r = \sum_k \sum_t \sum_n S_{n,k,t,r} w_{n,k,t,r} \tag{8}$$

Finally, if we are interested in an alternative measure of inflation (or a price index PI , which, as readers may have noted, does not require elementary prices to

be collected and recorded), we return to formula (3) and change it slightly to

$$PI_{k,n,t,r} = S_{k,n,t,r} = \frac{Y_{k,n,t,r}}{Y_{k,n,t=1,r}} \quad (9)$$

where an arbitrary year, normally the most remote, is used as a standard of reference. Formula (4) then becomes

$$PI_t = S_t = \sum_k \sum_n \sum_r S_{n,k,t,r} w_{n,k,t,r} \quad (10)$$

However, this does not solve all problems. Here, as in previous applications [3], we find that both alternative clustering criteria and, for any given criterion, different numbers of clusters result in slightly different values of our indices. As for the first problem (type of clustering), Ward appears to be the best criterion: its results are the most stable as the number of clusters vary. As for the second, there is an obvious trade-off between the number of clusters (the more the better, because the compared households then are, by our own definitions, more and more homogeneous), and the number of observations within clusters—where “observations” are households of all dimensions N , in all regions R , in all years T . It is however difficult to determine where the optimum lies.

4 An Empirical Application (Italy, 2004–2010)

This method has been recently introduced and applied to expenditure data for the estimation of equivalence scales in Italy in the years 2003–2008 [3]. In this paper we extend it in two ways: we estimate not only equivalence scales, but also *PPP* and inflation, and we apply it to two different datasets both the Istat yearly expenditure survey and the Bank of Italy *SHIW*—Survey of Household Income and Wealth, in the years 2004, 2006, 2008 and 2010. The *SHIW* survey is carried out every other year, on about 8000 households (24,000 individuals), in about 300 Italian municipalities. It is focused on household incomes and savings: there are a few questions on expenditure, including one on food expenditure, but they are not detailed enough, so that “standard” estimation methods (complete demand systems, based on expenditure shares) cannot be used in this case (see [13]).

For our method we need well-behaved indicators of economic well-being, and the four that we could compute from the *SHIW* survey are:

1. the food share (expenditure for food as a fraction of total expenditure),
2. a subjective assessment of the adequacy of the economic resources of the household. This derives from the six-scale answer to a question on “how difficult/easy is it for your family to make ends meet?” (possible answers are: very difficult, difficult, relatively difficult, relatively easy, easy, very easy),

3. the ratio between actual household consumption and a subjectively defined poverty line,
4. the saving share (the ratio between saving and income).

5 Main Results

Table 1 summarises our most important results. We have three different indices (equivalence scales, *PPP* and inflation), each with four outputs: two come from the Istat expenditure survey (with different sets of well-being indicators: see notes to the table), and two from the Bank of Italy *SHIW*. Of these, the former is based on household income and the other on household consumption, both reported by *SHIW* respondents.

As for the equivalence scale, the results from *SHIW* data confirm those found on expenditure data: large households cost more than small ones, but the increase is small, much smaller than Engel's food share would indicate (see, e.g., [6]), and comparable to that of the OECD [10] equivalence scale, obtained as the square root

Table 1 Main indices (equivalence scales, *PPP* and inflation) found for Italy, 2004–2010

	Household dimension (equivalence scale)				
	1	2	3	4	5
Income (<i>SHIW</i>)	1	1.473	1.807	2.023	2.101
Consumption (<i>SHIW</i>)	1	1.432	1.737	1.922	2.025
Consumption (Exp. Survey) ^a	1	1.282	1.542	1.675	1.781
Consumption (Exp. Survey) ^b	1	1.368	1.646	1.810	1.909
	Region (<i>PPP</i>)				
	NW	NE	Centre	South	Islands
Income (<i>SHIW</i>)	1	0.969	1.023	0.763	0.734
Consumption (<i>SHIW</i>)	1	0.947	1.015	0.754	0.742
Consumption (Exp. Survey) ^a	1	0.979	0.925	0.759	0.713
Consumption (Exp. Survey) ^a	1	1.003	0.934	0.745	0.707
	Year (inflation)				
	2004	2006	2008	2010	
Income (<i>SHIW</i>)	1	1.122	1.143	1.179	
Consumption (<i>SHIW</i>)	1	1.115	1.130	1.183	
Consumption (Exp. Survey) ^a	1	1.042	1.058	1.041	
Consumption (Exp. Survey) ^b	1	1.052	1.080	1.078	
<i>CPI (Istat)</i>	<i>1</i>	<i>1.038</i>	<i>1.089</i>	<i>1.115</i>	

See Fig. 1. All own estimates based on 100 clusters

CPI—Istat consumer price index (<http://www.istat.it/it/archivio/30440>)

^aWith dummy variables (=presence of durables). See De Santis and Maltagliati [3]

^bWithout dummies

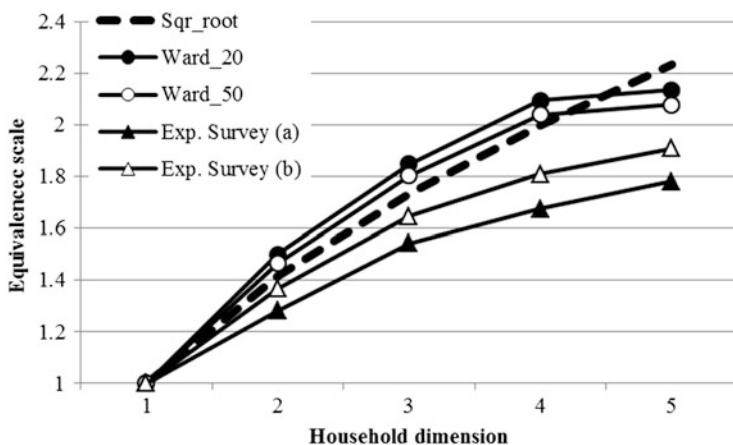


Fig. 1 Equivalence scales estimated for Italy, 2004–2010. *Notes:* Our estimates are obtained with the methodology described in this paper, using Ward as a clustering criterion, and total expenditure as a measure of resources. For income data (Bank of Italy SHIW), only the steepest (obtained with 20 clusters) and the flattest scales (with 50 clusters) are represented, highlighting how robust the estimates are to a change in the number of clusters. For the expenditure survey (Istat), the figures are those obtained with 100 clusters. The difference between the two versions, (a) and (b), depends on the choice of the indicators of well-being, as explained in the notes to Table 1. Households with more than five members have been discarded: the number of observations is 31,202 for SHIW and 187,482 for the expenditure survey (four waves in both cases). *Source:* own elaboration on Bank of Italy SHIW (Survey on Household Income and Wealth) data and Istat expenditure survey data, years 2004–2010

of household dimension (Fig. 1). SHIW data (in either of the two versions) yields an equivalence scale that is steeper, and therefore more convincing (or at least closer to prevalent expectations) than the one that comes from the expenditure survey.

When it comes to the estimation of PPP, or regional price differences, we basically once again obtain the same results with SHIW data (either scale) as with expenditure data. Prices appear to be more or less the same in the Centre-North, with minor differences within this area, but about 25 % lower in the South and Islands. These differences have never been measured precisely in Italy, but they are close to common perception, and are in line with the partial empirical estimates that have recently become available [2, 5].

There is however some inconsistency relating to the estimation of inflation. In this case there is an extra line: the official estimation of inflation in Italy (Istat CPI—Consumer Price Index). The estimates obtained here go qualitatively in the right direction (prices increase over time), but, quantitatively speaking, they do not fit well with official data. Using expenditure data, inflation is strongly underestimated if we form clusters using also dummy variables that indicate whether the households own certain durables (e.g. air conditioner; case “a” in Table 1). The reason for this is probably that in years of crisis (as in 2009 and 2010) households reduce their expenditure (which biases downward our estimates of inflation: see Eq. 9) but

they still possess the durables they had acquired previously, and therefore appear to be relatively well off. If durables are instead excluded from the procedure (i.e. not used as indicators of economic well-being) as in (b), our estimates of inflation get closer to the official value in their general trend, although not as much in each sub-period. In opposition to this, on the basis of the Bank of Italy *SHIW* inflation is overestimated: this is probably due the importance of subjective indicators (e.g. “can you make ends meet?”: the answer to this type of questions may reflect one’s concerns for the future, and vary with macroeconomic conditions and perspectives).

In both cases (Istat and Bank of Italy) these inconsistencies merit further analysis, but such results may be taken as a reminder that the standard of living is an elusive notion, with no clear empirical indicator (or set of indicators), and that all attempts at evaluating it can only give rough approximations at best. If the standard of living is measured incorrectly, so are all the measurements that derive from it, including the indices (equivalence scales, inflation and *PPP*) that this paper tries to estimate.

6 Discussion

In this paper we start from the idea (introduced in [3]) that equivalence scales can be estimated once we accept the assumption that the standard of living is a latent variable, correlated with observable (or manifest) variables, which, albeit imperfectly, indicate how rich or poor households are. Rules of consistency are defined for the selection of these “well-behaved” indicators of economic well-being and, when a sufficient number of them become available, they are used to form homogeneous clusters of households assumedly with the same standard of living. Once this is done, the estimation of cluster-specific equivalence scales is trivial (ratios between average incomes), and so is the following step of calculating an average of household specific scales, in order to estimate the general equivalence scale.

This paper introduces two novelties in this approach. The first is that the same logic may be used to estimate two more indices: inflation and *PPP*. Both of them are relevant and otherwise difficult to estimate, especially in the case of *PPP*. The second novelty lies in the use of alternative databases for the estimation of equivalence scales and other indices. The Bank of Italy *SHIW* (Survey of Household Income and Wealth) was used for the present application, but other datasets could be considered in the future: for instance, the LIS (Luxembourg Income Study) database, the ECHP (European Community Household Panel), the EU-SILC (European Survey on Income and living Conditions), etc. This extension is possible because of the flexibility of the definition that an indicator of economic well-being is virtually any variable that, everything else being equal, evolves consistently with resources.

Empirical and theoretical problems are still numerous: e.g. what indicators to select or drop when enough are available, how to form clusters, and how many of them to form. All of these choices affect the final results (estimates of equivalence scale, inflation or *PPP*), and this constitutes a cause for concern since no clear

indication has thus far emerged as to how to proceed in case of doubt. But a much higher uncertainty, even if normally hidden (with exceptions: see, e.g., [1]), also surrounds all the other estimates of equivalence scales, from the simplest (Engel's food share) to the most elaborate (complete demand systems). Several estimates of the same index can be produced with relative ease, which offers the additional advantage of providing a measure of how wide the margin of error is on the estimated variable of interest.

Acknowledgment Local funding from the University of Florence is gratefully acknowledged.

References

1. Bollino, C.A., Perali, F., Rossi, N.: Linear household technologies. *J. Appl. Econom.* **15**, 275–287 (2000)
2. Cannari, L., Iuzzolino, G.: Le differenze nel livello dei prezzi al consumo tra Nord e Sud. *Questioni di Economia e Finanza, Occasional Papers of the Bank of Italy*, 49, (2009)
3. De Santis, G., Maltagliati, M.: Clusters and equivalence scales. In: Torelli, N., Pesarin, F., Bar-Hen, A. (eds.) *Advances in Theoretical and Applied Statistics*. Springer, Berlin (2013)
4. Engel, E.: *Die Lebenskosten belgische Arbeiter-Familien früher und jetzt*. Heinrich, Dresden (1895)
5. Istat: *Le differenze nel livello dei prezzi al consumo tra i capoluoghi delle regioni italiane*. Rome (2010)
6. Istat: *La povertà in Italia nel 2011*. Rome. <http://www.istat.it/it/archivio/66983> (2012)
7. Kish, L.: *Survey sampling*. New York, Wiley, (1965[1995])
8. Lewbel, A., Pendakur, K.: *Equivalence Scales*. Entry for *The New Palgrave Dictionary of Economics*, 2nd edition, Boston College and Simon Fraser University. <http://www.sfu.ca/~pendakur/palequiv.pdf> (2006)
9. Muellbauer, J., van de Ven, J.: *Estimating equivalence scales for tax and benefits systems*. NIESR (Natl. Inst. of Econ. and Soc. Res.) discussion papers, 229. <http://www.niesr.ac.uk/pubs/dps/dp229.pdf> (2004)
10. OECD: *Growing unequal? Income distribution and poverty in OECD countries*. OECD, Paris (2008)
11. Perali, F.: *The Behavioral and Welfare Analysis of Consumption*. Kluwer Academic Publishers, Boston (2001)
12. Pollak, R.A., Wales, T.J.: *Welfare comparisons and equivalence scales*. *Am. Econ. Rev. Pap. Proc.* **69**, 216–221 (1979)

Websites

13. Bank of Italy *SHIW* <http://www.bancaditalia.it/statistiche/indcamp/bilfait>
14. ECHP project <http://circa.europa.eu/irc/dsis/echpanel/info/data/information.html>
15. EU-SILC project http://epp.eurostat.ec.europa.eu/portal/page/portal/microdata/eu_silc
16. LIS project <http://www.lisdatacenter.org/>

Learning-by-Exporting and Productivity: Evidences from a Panel of Manufacturing Firms

Maria Rosaria Ferrante, Marzia Freo, and Alessandro Viviani

Abstract

This paper investigates the dynamics of productivity experienced by firms that start to exporting. The effect of entering into international markets is disentangled by the self-selection component and its empirical distribution is evaluated. Results show that the impact on productivity of moving from the status of non-exporter to the status of exporter is different at different sections of the productivity distribution.

1 The “Learning-by-Exporting” Issue

A number of studies in the last decade, inspired by the seminal empirical work of Bernard and Jensen [1], have focused on the relationship between the internationalisation of firms and their productivity. In general, the main theme of these studies is that productivity is positively related to international involvement (for a review, see [2,3]).

This paper focuses on the debated issue of the learning-by-exporting (hereinafter LBE) hypothesis which states that firms that move from the status of non-exporter to the status of exporter (starters) experience an increase in productivity during the period following their entry into the export market. In fact, although an extensive stream of empirical literature on international trade shows that the most

M.R. Ferrante • M. Freo

Department of Statistical Sciences, University of Bologna, Bologna, Italy

e-mail: maria.ferrante@unibo.it; marzia.freo@unibo.it

A. Viviani (✉)

Department of Statistics, University of Florence, Florence, Italy

e-mail: viviani@disia.unifi.it

© Springer International Publishing Switzerland 2016

G. Alleva, A. Giommi (eds.), *Topics in Theoretical and Applied Statistics*,

Studies in Theoretical and Applied Statistics,

DOI 10.1007/978-3-319-27274-0_26

productive firms undergo a self-selection process to enter foreign markets, few empirical contributions support the LBE hypothesis. Besides, most of the literature focusing on the test of the LBE hypothesis estimates the effect by using either (1) regression models which estimate productivity premium by controlling for the relevant covariates or (2) methods developed in the context of the evaluation literature, such as the propensity score matching approach. Both of these approaches present strengths and weaknesses, but, primarily they limit the analysis to the estimation of the average treatment effect.

In this paper, we follow the strand of literature that purports that the international connections of firms have a significant association with firms' productivities and that this effect is different for the different points of the firms' performance distribution [4–12]. More precisely, we estimate the whole distribution of the net productivity premium caused by the entrance into international markets by disentangling the (raw) LBE premium from the component due to differences in the structural composition among three groups of firms: export starters firms, domestic firms and incumbent firms that entered international markets before the considered period and continue to export. In order to make inference on counterfactual distribution the Quantile Decomposition (QD) approach [13] is adopted. We obtain the estimate of the post-entry effect distribution by comparing the observed productivity distribution of the export starters group to the productivity distribution of its counterfactual.

The analysis is based on a newly available panel data set of firms recently developed by the Italian National Institute of Statistics [14]. These very rich longitudinal firm-level data allow for the solving of potential endogeneity problems and the proper evaluation of the causal effect of exporting on the performance of firms by also dealing with the self-selection bias.

2 Data and Methodology

A deep understanding of the structures of an economic system requires a profound knowledge of the time-dynamics concerning the most relevant economic phenomena. Obviously, the possibility to analyse the dynamics of economic system at the micro-level is constrained by data availability. With respect to a cross-sectional data set, panel data sets have significantly extended the research potential of data by allowing dynamic analyses and control for unobserved heterogeneity.

To meet the information need expressed by researchers, several European official statistical institutes have recently concentrated on designing a second generation of data sets derived from the integration of data already collected from surveys, census reports or administrative documents [15,16]. Obviously, this type of information is of a different nature with respect to the first generation of firm panel data sets, where a sample of firms or individuals was observed for multiple time periods.

In this perspective, during the last 10 years, ISTAT has promoted significant effort to build innovative longitudinal information on micro-data at the firm level, the Micro.3 database. Note that this type of information offers strong research potential for the dynamic analysis compared to the information collected in a survey. In fact,

the latter is generally limited to a few variables with a focus on one topic, whereas *Micro.3* combines the features of rich cross-sectional surveys with the potential for long observational periods [14]. Because of the importance of observing the behaviour of the same firm over successive time periods with references to variables from different fields (firm behaviour, financial data, structure and so on) in analysing the LBE effect on productivity, we base the research presented here on *Micro.3* data.

This database, which covers the period 1998–2007, contains information on firms with more than 20 employees. It is an integrated data system arising from the following three sources: (1) the census of Italian firms, that is, the SCI (*Sistema dei Conti delle Imprese*) database; (2) the PMI (*Piccole e Medie Imprese*) survey, which focuses on small and medium enterprises; and (3) the annual reports of incorporated firms collected by the Central Balance-Sheet Data Office. More technically, *Micro.3* is a catch-up panel, in which a cross-sectional data set is chosen at some time in the past and the units of analysis are then located in the present by subsequent observations. The validity of the database is largely supported by its census nature, which avoids possible biases in the data collection process, and its representativeness has been analysed in Biffignandi and Zeli [17].

2.1 Estimating Firms Productivity

To investigate the link between exporting and productivity with respect to performance outcomes, we consider Total Factor Productivity (hereinafter TFP) from an estimated two-factor Cobb–Douglas production function. In this framework, the main difficulty of estimating productivity at the micro-level arises from the simultaneity problem. In fact, the standard OLS estimates of the parameters of the log-transformed production function are biased as at least part of the TFP is known by the firm and influences the firm's input decisions such that the error terms and the regressors of the production function are correlated. The simultaneity problem between input decisions and productivity shock is accounted for by using the semi-parametric estimator proposed by Levinsohn and Petrin [18]. The TFP is measured at firm level by estimating eight Cobb–Douglas production functions by industry with value added as output, total costs of labour as labour input and the book value of fixed and intangible assets as capital input. All nominal variables are deflated by proper index numbers and deflated intermediate costs for goods and services are assumed as a proxy of capital (the results referring to the estimates are available from the authors). Because the level of TFP cannot be measured in any meaningful units, movements relative to a representative firm must be computed. To compute the relative movements, firm-specific TFPs are divided by the industry means. Note that these scaled TFPs are log transformed. Hereafter, when we refer to the estimated TFP, we always mean these log-transformations, which provide relative measures of how firm-specific TFPs diverge from the industry means during the period.

The raw TFP premium measures the percentage of the productivity differential between the TFP of exporter and non-exporter firms. From an unconditional perspective, exporter firms are found to be, at the median, from 9 to 13 % more

Table 1 Test on equality of median TFPs: export premiums by year

	Median premiums									
	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007
TFP	0.109	0.130	0.122	0.122	0.112	0.100	0.113	0.123	0.109	0.093
s.e.	0.011	0.011	0.010	0.011	0.012	0.013	0.012	0.014	0.016	0.015
Labour productivity	0.138	0.159	0.146	0.148	0.136	0.141	0.122	0.133	0.147	0.114
s.e.	0.012	0.013	0.015	0.015	0.016	0.016	0.014	0.014	0.018	0.019
Total output	0.759	0.823	0.773	0.765	0.778	0.786	0.743	0.817	0.773	0.637
s.e.	0.046	0.045	0.037	0.047	0.052	0.062	0.053	0.049	0.052	0.061
Total employment	0.569	0.597	0.606	0.564	0.631	0.571	0.612	0.540	0.482	0.396
s.e.	0.026	0.029	0.028	0.035	0.035	0.037	0.042	0.041	0.037	0.037
Capital stock	0.528	0.568	0.629	0.488	0.526	0.568	0.593	0.742	0.612	0.541
s.e.	0.059	0.055	0.057	0.064	0.056	0.065	0.063	0.058	0.050	0.066

productive than non-exporter firms in the years from 1998 to 2007. Exporters also employ 40–63 % more workers, produce more than 60 % more output and are strongly capitalised (Table 1). The differences reported are significant in that they show that exporters are very different from non-exporters both in TFP and covariates.

2.2 Estimating the Counterfactual Distribution

Starting from the raw observable internationalisation premium between two groups, Quantile Treatment Effects on Treated (QTET), which are the net internationalisation premiums at specified quantiles, are estimated by using the QD technique. QD allows to decompose the raw premium into a component explained by the differences in the pre-entry heterogeneity and a component explained by QTET. To properly estimate the net post-entry effect, it compares the observed productivity distribution of the starters to their counterfactual TFP distribution that is the distribution they would have had if they did not enter international markets. It intuitively works as follows. Groups are recognised different in many characteristics and effects of the treatment are assumed to differ across groups in the way they depend on these characteristics. Then starting from quantile conditional model of the TFP, Y , given the covariates X within the j -th group, $Q_{Y_j|X_j}(\tau|x)$, the conditional distribution implied for $Y_j|X_j$ is

$$F_{(Y_j|X_j)}(y|x) = \int_{(0,1)} I \{ Q_{Y_j|X_j}(\tau|x) \leq y \} d\tau$$

Given the covariate distribution $F_{X_i}(x)$, QD obtains estimates of functionals of the marginal distribution of the outcome estimated by integrating the conditional quantile model over the distribution of X . What it is interesting is that being the marginal distribution a function of the joint distribution of covariates and conditional distribution implied by quantile regressions, whose estimates are known after the estimation, it becomes possible to simulate counterfactual marginal distributions combining different scenarios for covariates and coefficients. For instance, it is possible to compute the counterfactual TFP distribution of firms of the group 1 provided that these firms had the same characteristics of firms of the group 2

$$F_{Y(1|2)}(y) = \int_{\mathcal{X}_2} F_{(Y_1|X_1)}(y|x) dF_{X_2}(x)$$

Finally, focusing on effects at θ unconditional quantile, the observed differences between the marginal TFP distributions over the groups are decomposed into a component which measures a QTET owing to LBE effect and a component explained by the differences in the composition of covariates owing to selection bias, as follows:

$$Q_{Y(2|2)}(\theta) - Q_{Y(1|1)}(\theta) = [Q_{Y(2|2)}(\theta) - Q_{Y(1|2)}(\theta)] + [Q_{Y(1|2)}(\theta) - Q_{Y(1|1)}(\theta)]$$

The QD has been run using the Stata program written by Chernozhukov et al. [13].

2.3 Building the Quasi-Experimental Design

To properly estimate the post-entry effect, we decompose the raw premiums in a portion due to the self-selection bias, caused by the pre-entry differences between groups, and a portion representing the net premium due to the entry of firms into international markets.

We focus on three groups of firms: incumbent exporters that are observed to export along the 10-year time window, domestic firms that are observed not to export during 9 out of the 10 years of the time window and a further group of firms that are observed to start exporting. Unfortunately, each year, a very small fraction of firms starts and continues exporting; thus, this last group has a small sample size. To capture as many observations as possible, we define starters as all firms that are observed to export for at least 5 years after having not exported for the 2 previous years. For example, consider that for each firm we observe a 10-year sequence of export dummies. Firms are considered incumbents if they show a full one sequence (1-1-1-1-1-1-1-1-1-1), domestics if within the sequence only one 1 may be retrieved and starters if their sequence nests a sub-sequence of this type 0-0-1-1-1-1-1-1. Then, because starter firms may have begun to export in a year between 2000 and 2003,

Table 2 Number of TFP observations by status and time to treatment

Time to treatment (j)	-2	-1	+1	+2	+3	+4
Incumbents	1799	1811	1783	1787	1782	1779
Starters	197	205	196	195	192	193
Domestics	351	346	262	207	169	154

Table 3 Averages and per cent distribution of key variables at time $t = -2$ by status

	Incumbents	Starters	Domestics
TFP (raw premium vs. domestics)	0.279	0.159	–
Per cent of producers of intermediates goods	45.0	43.9	48.4
Durable goods and instruments	33.8	32.2	20.9
Non-durable goods	21.2	23.9	30.7
Per cent of producers located in North-West	49.5	41.1	23.8
North-East	36.9	33.8	20.4
Centre	9.7	15.5	19.7
South	3.9	9.6	36.1
Total output (thousand EUR)	76,244	50,206	18,754
Total employment	324	203	109
Capital stock (thousand EUR)	16,159	12,036	6643

we fix the treatment time $j = 0$ in the year when each firm enters the export market, and we align all observations with respect to the treatment time. Thus, we consider the time to treatment variable in terms of the advance or delay to the treatment time (since $j = -2$ to $j = 4$). After removing observations without information on the TFP, we have three groups whose number of observations is shown in Table 2. Note that the definition of starters group leads us to observe a 2-year long period before the treatment ($j = -2, -1$) and a 4-year long period after the treatment ($j = 1, 2, 3, 4$).

Firm heterogeneity can be found in the pre-entry compositional characteristics of the three groups of firms and, more importantly, in the pre-entry TFP that we observe at time to treatment $j = -2$ (Table 3). Within the domestics, more than one-half of the firms produce intermediate goods, and more than one-third of the firms (36.1 %) are located in the southern region of the country. The profiles by activity and location of starter and incumbent firms are quite similar. Incumbent and starter firms are strongly larger than domestic firms. For example, at a mean level, incumbent and starter firms have double the total employment of domestic firms, they have just less than double the capital stock, and they have three times the output.

3 Self-Selection and LBE: Empirical Results

To test the presence of self-selection in the observable pre-entry outcome, we investigate if the three previously defined groups differ significantly in TFP with a time reference that is before the treatment time $j = 0$. More specifically, we compare

Table 4 Raw TFP premium: pre-entry levels and dynamics

	Incumbents vs. domestics			Starters vs. domestics			Starters vs. incumbents		
Quantile	0.2	0.5	0.8	0.2	0.5	0.8	0.2	0.5	0.8
Level	0.312**	0.264**	0.300**	0.171**	0.159**	0.206**	-0.140**	-0.106**	-0.094°
Growth	0.044°	0.011	-0.003	0.010	0.000	-0.009	-0.034	-0.011	-0.006

°, ** significant at 90 and 99 % (Bootstrap confidence intervals)

the three pairs of groups: incumbents vs. domestics, starters vs. domestics and starters vs. incumbents. We focus on all the TFP level and TFP growth rate to verify if the self-selection is on the pre-entry levels and/or on the trends. To this end, first we test whether differences across groups in the TFP level at time to treatment $j = -2$ are significant and then whether the dynamics of the TFP is significantly different across groups from time to treatment $j = -2$ to $j = -1$. The observed raw premiums before-entry (measured by differences in TFP level across groups) appear significantly different from zero for both the realised comparisons (incumbents vs. domestics, starters vs. domestics and starters vs. incumbents). At time to treatment $j = -2$ incumbent firms are observed (in median) to be 26.4 % more productive than domestic firms, while starter firms are 15.9 % more productive than domestics and 10.6 % less than incumbents (Table 4). Looking at the distributions, the differences are higher at the extremes, that is, for the highest and lowest TFP values. Lower performing incumbents and starters are, respectively, 31.2 and 17.1 % more productive than less performing domestic firms (at quantile 0.20), and the best performing incumbents and starters are 30.0 and 20.6 %, respectively, more productive than the best performing domestic firms (at quantile 0.80). Thus, the possible presence of different dynamics before entry is considered. If present, this must be accounted for to appropriately evaluate the post-entry effect. The findings indicate significant positive dynamics (even if significant at 90 %) in the raw premium of less performing incumbents with respect to less performing domestics. All the other premiums' dynamics are not significantly different across groups before the entry period (Table 4). Because of the scant significance and to improve the homogeneity and readability of the results, we assume here the absence of different dynamics.

We estimate the net TFP premium obtained by removing the selection bias from the raw premium. More specifically, from the perspective of the decomposition approach, the difference of the TFP across groups measures the raw export premium, and it can be decomposed into two parts: the component owing to characteristics that may be interpreted as the effect of self-selection on observables and the component owing to coefficients that may be interpreted as the estimate of the treatment effect. To investigate the performance of starter firms in the post-entry period, we pose the treatment period at $j = 0$ and focus on four post-entry periods denoted by $j = 1, 2, 3, 4$.

The first step of the QD approach consists of estimating quantile regression models and explaining raw TFP for the different groups of starters (s), domestics (d) and incumbents (i). With this aim, we specify the following τ -th quantile regressions

model that we estimate over each group ($g = s, d, i$) and different time to treatment $j = 1, 2, 3, 4$:

$$TFP_{igj} = \alpha^g(\tau) + \sum_{k=1}^2 \beta_k^g(\tau) IND_{ig,k} + \sum_{l=1}^3 \gamma_l^g(\tau) REG_{ig,l} + \eta^g(\tau) Z_{ig,-2} + u_{igj}(\tau)$$

where TFP_{igj} is the TFP level at time to treatment j of firm i , which belongs to group g ; $IND_{ig,k}$ are two out of three dummies indicating the principal industry group to which firm i of the group g belongs at year 1999; $REG_{ig,l}$ are three out of four dummies indicating the macro-area in which the firm i of the group g is located at year 1999; $Z_{ig,-2}$ is the vector of observable covariates that controls for selection bias, that is the pre-entry raw productivity level ($TFP_{ig,-2}$ for firm i at time to treatment $j = -2$).

The second step consists of estimating the net TFP through the QD approach, that is, of removing the selection bias due to different group compositions observed before entry and to raw TFP pre-entry levels. Results of the decomposition are reported in Table 5. The estimated net TFP premiums of starters versus domestic

Table 5 Decomposition of the TFP level

TFP <i>quantile</i>	Raw differences			Net differences		
	0.2	0.5	0.8	0.2	0.5	0.8
<i>Time to treatment (j)</i>						
<i>Incumbents vs. domestics</i>						
-1	0.315**	0.289**	0.292**	0.077**	0.069**	0.071*
0	0.364**	0.318**	0.334**	0.132**	0.124**	0.140**
1	0.366**	0.281**	0.261**	0.123**	0.079*	0.059
2	0.353**	0.264**	0.209**	0.126**	0.062	0.017
3	0.212**	0.189**	0.129**	0.102	0.039	-0.039
4	0.238**	0.215**	0.155**	0.086	0.054	0.037
<i>Starters vs. domestics</i>						
-1	0.155**	0.161**	0.177**	0.037	0.036	0.018
0	0.237**	0.193**	0.221**	0.100**	0.073*	0.081
1	0.199**	0.146**	0.153**	0.067 ^o	0.029	0.014
2	0.205**	0.157**	0.115*	0.086*	0.043	-0.007
3	0.102*	0.084*	0.047	0.051	0.004	-0.053
4	0.160**	0.130**	0.055	0.078	0.038	-0.011
<i>Starters vs. incumbents</i>						
-1	-0.160**	-0.128**	-0.115**	-0.032	-0.030*	-0.052
0	-0.127**	-0.126**	-0.113**	-0.020	-0.037 ^o	-0.054
1	-0.167**	-0.135**	-0.108*	-0.075**	-0.055**	-0.056
2	-0.149**	-0.106**	-0.094*	-0.070*	-0.037 ^o	-0.048
3	-0.111**	-0.106**	-0.081 ^o	-0.034	-0.039 ^o	-0.037
4	-0.078*	-0.085**	-0.100**	-0.004	-0.018	-0.058

^o, **, * significant at 90, 95 and 99 % (Bootstrap confidence intervals)

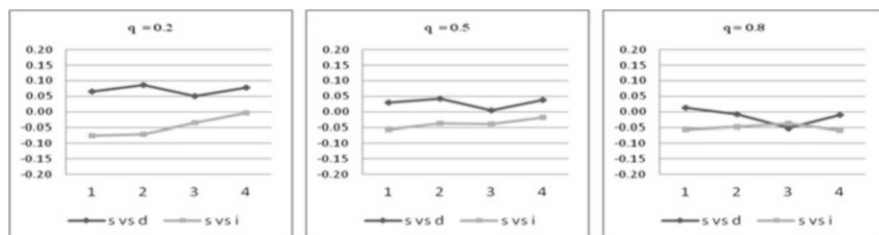


Fig. 1 Net TFP premiums at time to treatment $j = 1, 2, 3, 4$: starters vs. domestics (s vs d) and starters vs. incumbents (s vs i) at quantiles (q) 0.2, 0.5 and 0.8

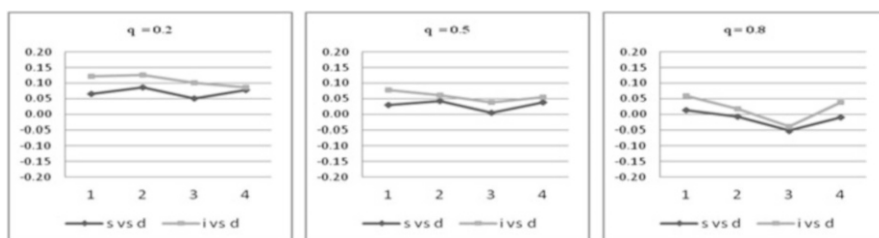


Fig. 2 Net TFP premiums at time to treatment $j = 1, 2, 3, 4$: starters vs. domestics (s vs d) and incumbents vs. domestics (i vs d) at quantiles (q) 0.2, 0.5 and 0.8

firms are positive over the post-entry period in the first half of the distribution. At the same time, the net premiums of the starter firms in the higher half of the TFP distribution are null or negative (Fig. 1). This suggests that, while TFPs of lower and medium performing starter firms overcome TFPs of lower and medium performing domestic firms, this does not occur for better performing firms. With regards to the level of the estimated net TFP premiums of starters versus incumbents, they are negative for low, medium and high performing firms. The evolution of the TFP premiums during the post-entry periods may be usefully captured by jointly considering empirical evidences from Figs. 1 and 2 in two directions. On the one hand, net premiums of starter with respect to domestic firms along the post-entry period are quite stable, the only exception being the fall in period $j = +3$. On the other, the dynamics of net premiums of starters with respect to incumbent firms is positive. Thus, the negative gap of starters tends to decline and becomes null either at time to treatment $j = +3$ for the comparison that entails lower performing firms, at time to treatment $j = +4$ for the comparison that entails medium performing firms or suddenly at time to treatment $j = +2$ for the comparison that entails higher performing firms. A first finding is that starters, in the post-entry periods, increase their TFPs more than incumbents but do not increase their TFPs more than domestics.

That is, the entry into international markets produces an acceleration of starter TFPs with respect to incumbents. To better understand this finding, it must be remembered that the analysed years represent a phase of a negative cycle of Italian

exports during which premiums of incumbents diminished as long as they had become not significant (Table 5). According to expectations, if a post-entry effect had not intervened, the net premiums of starters would have experienced a decrease of the same order of incumbents' premiums. On the contrary, the levels of premiums for starters remained quite stable during the post-entry period and converged to the premiums of incumbents.

At time to treatment $j = +4$, premiums of lower and medium starter firms may not be distinguished from the ones of lower and medium incumbent firms (Fig. 2). With regards to the firms in the top section of the TFP distributions, the net premiums are usually lower than in the other sections of the distributions, and they are never significant, indicating that the TFPs of the best performing firms are not sensitive to internationalisation status. The main findings draw important conclusions in terms of economic policy. In fact, if a learning-by-exporting effect is present, policies should remove obstacles to export entry as this may help to increase the number of firms that successfully act on the world market in the future and that contribute to economic growth through the increase of their productivity; on the contrary, if the learning-by-exporting effect is absent, policies should focus on directly fostering productivity.

4 Concluding Remarks

Many papers in the empirical literature on the heterogeneity of firms find that a productivity premium is associated with international involvement. Two fundamental questions from this topic are still unanswered. The first question involves the amount of the productivity premium, which, although in relative terms, is rarely estimated. The second question involves the direction of the causality link between productivity and international openness. Does a self-selection mechanism induce the more productive firms to enter the international market, or do internationalised firms, under the pressure of global competition, become more productive by means of a learning-by-exporting process?

The present paper investigates these issues during a negative cycle for Italian exporter performances. At first glance, the raw TFP premiums for the exporting firms are found to be quite high according to many firms' characteristics, such as TFP, labour productivity, total output, capital stock and total employment.

TFP estimates are next decomposed with a QD approach. Through this technique, the overall productivity gaps between the three groups of firms, incumbents, starter exporters and domestics are disentangled, separating the part of the gap that is explained by differences in firm characteristics, that is, the different pre-entry composition of the groups, from the part that is owing to internationalisation.

The main findings support the intervention of an LBE effect in the post-entry period for firms that begin exporting, even if not uniformly distributed across firms. In the present study, the LBE is for the lower and medium performing starter firms. Firms in the bottom half of the TFP distribution experience an increase with respect to incumbent firms but not compared to domestic firms after entering the

international market. Thus, internationalisation has resulted in at least one-half of the firms improving their performance relative to incumbents.

Net premiums of the best performing firms from all three groups are statistically not significant, indicating that firms that perform at their best may operate both in domestic and export markets. For these best performing firms, the policy indication should address the increase of TFP itself rather than the internationalisation.

Further work will involve refining the analysis. Details regarding year and cyclical effects should be more clearly disentangled, and the analysis should be extended to include additional performance features. Addressing these challenges should provide more robust results.

References

1. Bernard, A.B., Jensen, J.: Exporters, jobs and wages in US manufacturing: 1976-87. *Brook. Pap. Econ. Act.: Microeconomics* **1995**, 67–112 (1995)
2. Greenaway, D., Kneller, R.: Exporting and productivity in the United Kingdom. *Oxf. Rev. Econ. Policy* **20**(3), 358–371 (2004)
3. Wagner, J.: Export and productivity: a survey of evidence from firm-level data. *World Econ.* **30**(1), 60–82 (2007)
4. Bellone, F., Guillou, S., Nesta L.: To what extent innovation accounts for firm export premia? Technical report (University of Nice - Sophia Antipolis) (2010)
5. Dimelis, S., Louri, H.: Foreign ownership and production efficiency: a quantile regression analysis. *Oxf. Econ. Pap.* **54**(3), 449–469 (2002)
6. Falzoni, A.M., Grasseni M.: Home country effects of investing abroad: evidence from quantile regressions, CESPRI Working Paper, 170. pp. 1–37, Milano (2005)
7. Ferrante, M.R., Freo, M.: The total factor productivity gap between internationalised and domestic firms: net premium or heterogeneity effect? *World Econ.* **35**(9), 1186–1214 (2012)
8. Haller, S.: Intra-firm trade, exporting, importing and firm performance. *Can. J. Econ.* **45**(4), 1397–1430 (2012)
9. Powell, D., Wagner, J.: The exporter productivity premium along the productivity distribution: evidence from quantile regression with nonadditive firm fixed effects. *Rev. World Econ.* **150**(4), 763–785 (2014)
10. Serti, F., Tomasi, C.: Self-selection and post-entry effects of exports: evidence from Italian manufacturing firms. *Rev. World Econ.* **144**(4), 660–694 (2008)
11. Wagner J.: Exports and firm characteristics in German manufacturing industries. *Appl. Econ. Q.* **57**(2), 107–143, 145–160 (2011)
12. Trofimenko, N.: Learning by exporting: does it matter where one learns? Evidence from Colombian manufacturing firms. *Econ. Dev. Cult. Chang.* **56**, 871–894 (2008)
13. Chernozhukov, V., Fernández-Val, I., Melly, B.: Inference on counterfactual distributions. *Econometrica* **81**(6), 2205–2268 (2013)
14. Grazzi, M., Sanzo, R., Secchi, A., Zeli, A.: The building process of a new integrated system of business micro-data 1989–2004. *J. Econ. Soc. Meas.* **38**(4) 291–324 (2013)
15. Konold, M.: New possibilities for economic research through integration of establishment-level panel data of German official statistics. *Schmollers Jahrbuch/J. Appl. Soc. Sci. Stud.* **127**(2), 321–334 (2007)
16. Wagner, J.: The post-entry performance of cohorts of export starters in German manufacturing industries. *Int. J. Econ. Bus.* **19**(2), 169–193 (2012)

17. Biffignandi S., Zeli A.: Integrating databases over time: what about representativeness in longitudinal integrated and panel data? Paper presented at the European Conference on Quality in Official Statistics, Helsinki, Finland (2010)
18. Levinsohn, J., Petrin, A.: Estimating production functions using inputs to control for unobservables. *Rev. Econ. Stud.* **70**(2), 317–342 (2003)

A Multivariate VEC-BEKK Model for Portfolio Selection

Andrea Pierini and Alessia Naccarato

Abstract

The use of bivariate vector error correction models and Baba–Engl–Kraft–Kroner models is proposed for the selection of a stock portfolio (Markowitz portfolio) based on estimates of average returns on shares and the volatility of share prices. The model put forward is applied to a series of data regarding the prices of 150 shares traded on the Italian stock market (BIT) between 1 January 1975 and 31 August 2011.

1 Introduction

The selection of a stock portfolio is broadly discussed in the literature, generally with reference to heteroskedastic regression models [1]. The models used in the case of multiple time series are of the vector autoregressive (VAR) type [4].

This paper proposes the use of vector error correction (VEC) and Baba–Engl–Kraft–Kroner (BEKK) models for the selection of a stock portfolio. In other words, it addresses the problem of estimating average returns and the associated risk on the basis of the prices of a certain number of shares over time. This estimate is then used to identify the assets offering the best performance and hence constituting the best investments. While Campbell [4] proposes the use of a VAR (1) model, it is suggested here that use should be made of VEC models, which make it possible to take into account any cointegration between the series employed and the market trend as measured by means of the Thomson Reuters Datastream Global Equity Italy Index [5].

A. Pierini (✉) • A. Naccarato

Department of Economics, University of Roma Tre, Via S. D'Amico 77, 00145 Roma, Italy
e-mail: andrea.pierini@uniroma3.it; alessia.naccarato@uniroma3.it

Moreover, while Bollerslev et al. [2] employ diagonal vectorization (DVEC) models to estimate share volatility, the use of a BEKK model, as proposed here, makes it possible to extend the estimation procedure based on DVEC models so as to take into account also the correlation between the volatility of the series and the volatility of the market trend.

The series considered regard the Italian stock market (BIT), and specifically the monthly figures for the top 150 shares in terms of capitalization, from 1 January 1975 to 31 August 2011. The estimation procedure proposed for portfolio selection involves two phases.

In the first, a two-dimensional VEC model is developed for all of the 150 shares considered in order to obtain an estimate of the average stock market return. A BEKK model is then applied to the series of residuals thus obtained in order to estimate the volatility of the series.

The second regards the selection of shares for inclusion in the portfolio. Only those identified as presenting positive average returns during the first phase are considered eligible. For the purpose of selecting the most suitable of these, a new endogenous variable is constructed as the product of two further elements, namely the price-to-earnings ratio (P/E) and earnings per share (EPS). This variable, which indicates the “intrinsic value” of the share in question, is not constructed for the entire set of 150 shares but only for those presenting positive average returns in the first phase, as it would be pointless in the case of negative returns. The VEC-BEKK model is applied once again to this new series in order to estimate the intrinsic value of the shares, and the top 10, as suggested in [7], are selected for inclusion in the portfolio on the basis of the difference between this intrinsic value and the price estimated in the first phase.

A quadratic programming model is then employed to determine the quantities to be bought of each of the ten shares selected.

It should be noted that the variable $P/E \cdot \text{EPS}$ is estimated for each industrial sector, as suggested in [9].

2 Model Summary

A concise outline is now given of the phases involved in the selection of shares for inclusion in the portfolio as well as the quantity of shares to be bought for each type selected. The starting point is the $K = 150$ series, regarding the average returns $R_{k,t}$ on the shares, and the average return of the market $R_{M,t}$, $t = t_k, \dots, T$, $k = 1, \dots, K$. It should be noted in this connection that the length of the series considered is not homogeneous because not all of the joint-stock companies are quoted as from the same point in time. This aspect involves further complications in the estimation procedure.

Phase one: For each series, the model $VAR_2(p)$ is constructed for the random vector $y_t = [y_{1,t}, y_{2,t}]' = [R_{k,t}, R_{M,t}]'$

$$y_t = \mu_t + A_1 y_{t-1} + A_2 y_{t-2} + \dots + A_p y_t + u_t \tag{1}$$

with $\mu_t = \mu_0 + \mu_1 t$, A_i matrix 2×2 , $i = 1, \dots, p$ of the unknown coefficients, and $u_t = [u_{1,t}, u_{2,t}]'$ the vector of errors such that $u_t \sim N(0, \Sigma_u)$.

Model (1) can be rewritten as follows to take into account and possible cointegration of the variables considered:

This model can be rewritten in the form of $VEC_2(p-1)$, which shows manifestly the possible cointegration, which we use in presence of it

$$\Delta y_t = \mu_t + \Pi y_{t-1} + \Gamma_1 \Delta y_{t-1} + \Gamma_2 \Delta y_{t-2} + \dots + \Gamma_{p-1} \Delta y_{t-p+1} + u_t \tag{2}$$

The AIC criterion is used to estimate the lag \hat{p} , with reference to model (2), and the LR test is carried out to ascertain the presence of cointegration.

Finally, the method proposed by Johansen [12] is applied to obtain the maximum-likelihood estimation (MLE) of the parameters $\mu_0, \mu_1, \Pi, \Gamma_1, \dots, \Gamma_{p-1}$.

The Portmanteau test is used to ascertain the presence of correlation of residuals, the generalized Lomnicki–Jarque–Bera test for the normality of residuals, and the ARCH test to determine heteroskedasticity.

In the event of the latter test revealing the presence of heteroskedasticity, the BEKK(1,1) model [6] is used to estimate the conditional variance–covariance matrix $\Sigma_t = cov(u_t | \text{past}) = ((\sigma_{i,j}(t))_{i,j=1,\dots,n})$, which has the following structure:

$$\begin{aligned} \begin{bmatrix} \sigma_{11,t} & \sigma_{12,t} \\ \sigma_{21,t} & \sigma_{22,t} \end{bmatrix} &= \begin{bmatrix} a_{11} & 0 \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} a_{11} & a_{21} \\ 0 & a_{22} \end{bmatrix} \\ &+ \begin{bmatrix} a_{11,1} & a_{12,1} \\ a_{21,1} & a_{22,1} \end{bmatrix} \begin{bmatrix} u_{1,t-1}^2 & u_{1,t-1} u_{2,t-1} \\ u_{2,t-1} u_{1,t-1} & u_{2,t-1}^2 \end{bmatrix} \begin{bmatrix} a_{11,1} & a_{21,1} \\ a_{12,1} & a_{22,1} \end{bmatrix} \\ &+ \begin{bmatrix} b_{11,1} & b_{12,1} \\ b_{21,1} & b_{22,1} \end{bmatrix} \begin{bmatrix} \sigma_{11,t-1} & \sigma_{12,t-1} \\ \sigma_{21,t-1} & \sigma_{22,t-1} \end{bmatrix} \begin{bmatrix} b_{11,1} & b_{21,1} \\ b_{12,1} & b_{22,1} \end{bmatrix} \end{aligned} \tag{3}$$

Phase two: The estimates obtained in phase one are used to select the shares for which positive average returns are predicted. For the shares thus selected and for each industrial sector (IS), the model $VEC_2(p-1)$ –BEKK(1, 1) is estimated for the random vector $y_t = [y_{1,t}, y_{2,t}]' = [((P/E) \cdot (EPS))_{h,t}, ((P/E) \cdot (EPS))_{IS_h,t}]'$, where $h = 1, \dots, H$ is the index that identifies only the series with positive returns selected out of the initial 150.

On the basis of the $((P/E) \cdot (EPS))_{IS_h, T+1}$ and $R_{h, T+1}$ forecasts obtained in phase two, the shares are listed for each industrial sector in decreasing order with respect to the values of the difference between intrinsic value and expected price. The first $n = 10$ shares are thus selected to make up the portfolio.

Finally, in order to determine the quantities to be bought of each of the ten shares selected, it is necessary to solve the Markowitz problem [13] by estimating the matrix of share volatility. To this end, let be \hat{V}_t the estimator of the matrix $n \times n$ of volatility V_t for $t = T + 1$, the elements of which are $v_{i,j}(t) = \text{cov}(R_t|\text{past}), i, j = 1, \dots, n$. The elements of \hat{V}_t are given by:

$$\hat{v}_{i,j}(T + 1) = \begin{cases} \hat{\sigma}_{11,T+1|T}^{(i)} & \text{se } i = j \\ \hat{c}_{ij} & \text{se } i \neq j \end{cases} \tag{4}$$

with $\hat{C} = (\hat{c}_{i,j})_{i,j=1,\dots,n} = \sum_{t=t_{\max}}^T (R_{i,t} - \bar{R}_i)(R_{j,t} - \bar{R}_j)' / (T - t_{\max})$, $t_{\max} = \max\{t_i, t_j\}$, $i, j = 1, \dots, n, n = 10$. On the basis of (4), the solution of the quadratic Markowitz type problem

$$\begin{cases} \min \omega' \hat{V} \omega \\ \omega \geq 0 \\ \omega' \mathbf{1} = 1 \end{cases} \tag{5}$$

for the future time $T + 1$ can be obtained with the approximation given by the Goldfarb–Idnani dual method [8] that we briefly describe in the following.

The equality constraints in (5) can be seen as disequality constraints:

$$\omega' \mathbf{1} = 1 \Leftrightarrow \begin{cases} \omega' \mathbf{1} \geq 1 \\ \omega' (-\mathbf{1}) \geq -1 \end{cases}$$

So the problem (5) can be rewritten as:

$$\begin{cases} \min \omega' \hat{V} \omega \\ C' \omega - b \geq 0 \end{cases} \tag{6}$$

where

$$C' = \begin{bmatrix} 1 & \dots & 0 \\ 0 & \ddots & 1 \\ 1 & \dots & 1 \\ -1 & \dots & -1 \end{bmatrix}, b = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ -1 \end{bmatrix} \tag{7}$$

We call active set a subset of the $m = n + 2$ constraints in (6) that are satisfied as equalities by the current estimate ω_0 of the solution of (6).

A subproblem $P(J)$ of the problem (6) has the same objective function but only a subset of constraints indexed by $J \subset K = \{1, 2, \dots, m\}$.

Moreover if the solution ω_0 of a subproblem $P(J)$ lies on some linearly independent active set of constraints indexed by $A \subset J$ we call the pair (ω_0, A) solution pair.

The basic step of the dual algorithm are:

- Step 0: start with the solution pair $(\omega_0, A) = (\hat{V}^{-1}\mathbf{1}/(\mathbf{1}'\hat{V}^{-1}\mathbf{1}), \emptyset)$*
- Step 1: repeat until all constrains are satisfied:*
 - (a) choose a violated constraint $p \in K - A$*
 - (b) if $P(A \cup \{p\})$ is infeasible \Rightarrow stop the problem is infeasible*
 - (c) else obtain a new solution pair (ω_1, B)*
 - with $A_1 \subseteq A, B = A_1 \cup \{p\}$*
 - so that $\omega_1'\hat{V}\omega_1 \geq \omega_0'\hat{V}\omega_0$,*
 - and set $(\omega_0, A) = (\omega_1, B)$*
 - This is always possible by solving $P(A \cup \{p\})$.*
- Step 2: when all constrains are satisfied stop $\Rightarrow \omega_0$ is the solution.*

To obtain a better diversification we have also find the solution of the quadratic problem of Markovitz type (5) without the constraint $\omega \geq 0$, for the future time $T + 1$ using the explicit solution

$$\hat{\omega}_{\text{opt},T+1} = \hat{V}_{T+1}^{-1}\mathbf{1}/(\mathbf{1}'\hat{V}_{T+1}^{-1}\mathbf{1}) \tag{9}$$

Then we put to zero the shorting and repropionate the remaining $\omega_i, i = 1, \dots, n$. We omit the constraint of a fixed value for the expected return to eliminate the sensitiveness of allocation optimization to errors in predicted returns [11].

When the matrix \hat{V} is not positive definite, we propose the approximation with the nearest matrix in the Frobenius sense, retaining the same diagonal of \hat{V} with estimated elements given by the BEKK model part.

To find the nearest matrix H to \hat{V} we proceed as follows:

Firstly we find the matrix \hat{V}^c where $\hat{v}_{ij}^c = \frac{\hat{v}_{ij}}{\sqrt{\hat{v}_{ii}\hat{v}_{jj}}}$.

Then with the Higham algorithm [10] we find the nearest symmetric positive definite matrix with unit diagonal H^c to \hat{V}^c .

Finally we find the matrix H where $h_{ij} = h_{ij}^c \sqrt{\hat{v}_{ii}\hat{v}_{jj}}$.

To find H^c we solve the following problem

$$\begin{cases} \min_X \| \hat{V}^c - X \| \\ X \in S = \{Y = Y^T \in \mathfrak{R}^{n \times n} : Y_{\text{pos.def.}}\} \\ X \in U = \{Y = Y^T \in \mathfrak{R}^{n \times n} : y_{i,i} = 1\} \end{cases} \tag{10}$$

where $\| A \| = \sqrt{\sum_{i,j} a_{ij}^2}$ is the Frobenius norm.

The solution of (7) is found by iteratively projecting onto the subspaces S, U , after applying the Dykstra's correction [3] to the S projection to guarantee the

convergence. So the following algorithm is used:

$$\begin{aligned}
 \Delta S_0 &= 0, Y_0 = \hat{V}^c \\
 k &\in \{1, 2, \dots\} \\
 R_k &= Y_{k-1} - \Delta S_{k-1} \\
 R_k &= QDQ^T, D = \text{diag}(\lambda_i), Q = \{q_1, \dots, q_n\} \\
 &\text{where } \lambda_i \text{ eigenvalue of } R_k, \\
 &q_i \text{ the corresponding eigenvector, } i = 1, \dots, n \\
 \lambda^+ &= \text{diag}(\max(\lambda_i), 0) \\
 P_S(R_k) &= Q\lambda^+Q^T \\
 X_k &= P_S(R_k) \\
 \Delta S_k &= X_k - R_k \\
 \theta &= \text{diag}(X_k - I) \\
 P_U(X_k) &= X_k - \text{diag}(\theta_i) \\
 Y_k &= P_U(X_k) \\
 &\text{if } \|X_k - Y_k\| / \|Y_k\| \leq \epsilon \Rightarrow \text{stop} \\
 &\text{next}
 \end{aligned} \tag{11}$$

Boyle–Dykstra [3] show that X_k, Y_k converge to \hat{V}^c . However even if the repeated projections converge to the point H^c in the intersection between S and U nearest to the starting point \hat{V}^c , we cannot guarantee that the norm of the difference in (10) will be small.

Sometimes it may be better to relax the constraints in order to obtain a better approximation in terms of smaller norm as it is shown in the results.

3 Results

Application of the model proposed in this work to the monthly figures for the 150 BIT shares with the highest level of capitalization indicates the following results:

(a) An optimal lag p of 2–9 months, see Fig. 1a.

In particular, the optimal lag is 2 months for 77% of the entire set of 150 shares. This means that just 2 months of observation are sufficient to predict the average returns on the vast majority of the shares considered.

(b) The degree of cointegration proves equal to 2 for 84% (vin) or 91% (return) of the 150 shares, 1 for 12% (vin) or for 7% (return) and 0 for the remaining 4% (vin) or 2% (return), see Fig. 1b.

So there is the presence of cointegration in phase one and two indicating the need of VEC model to better explain the time series.

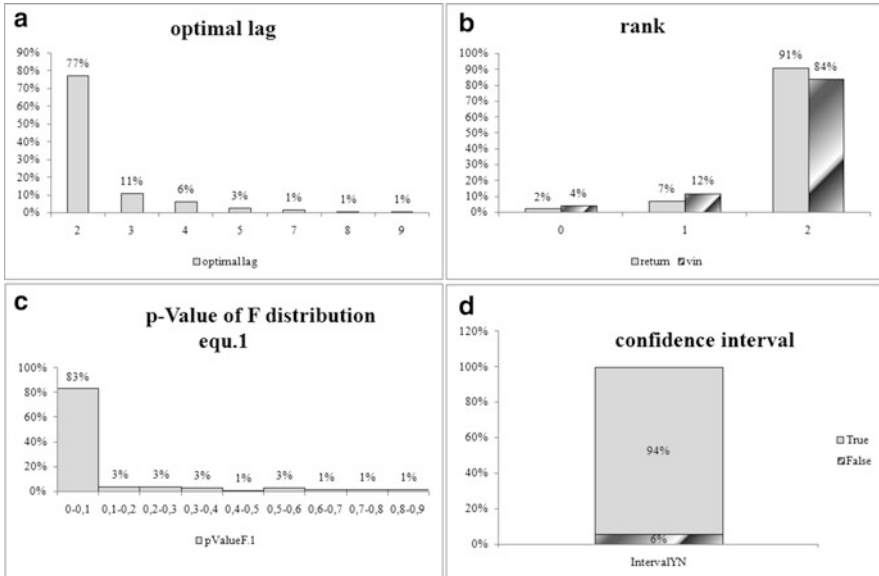


Fig. 1 Top left: distribution of optimal lag. Top right: distribution of the ranks. Bottom left: distribution of the *p* value of *F*. Bottom right: distribution of the confidence interval

- (c) The *coefficients* of the models estimated in both steps of the procedure prove significant for almost all of the series considered, that is to say for 86 % of the 150 shares the *p* value of the *F* statistic is less than 0.2 , see Fig. 1c.
- (d) The BEKK estimate of *volatility* for each share is between 0.001 and 0.01 for 93 % of the series and never above 0.021.
- (e) The *confidence interval* at the level of significance of 95 % contains the actual value_{*T*+1} in 94 % of the series, see Fig. 1d.

The VEC-BEKK model can therefore be considered reliable for most of the series for the purposes of prediction.

- (f) The best portfolio has a *monthly average return* of 1.9 %, a *monthly standard deviation* of 0.655, and a *Sharpe index* of 0.029.

It is obtained by ranking in decreasing order inside each of the ten industrial sector and selecting the first stocks. We call this ranking procedure *partial ranking*.

We also obtained a portfolio by ranking in decreasing order selecting the first ten stocks. We call this ranking procedure *total ranking*.

In Table 1 we see that different ranking procedure gives different portfolio.

The optimal partial ranking is the best portfolio because it has the greatest return and the lowest risk (volatility) and so the highest Sharpe. However the proportional total ranking is the best diversified portfolio because it has the greater number of nonzero stocks. It has a return near to the best one but an almost double risk.

Table 1 Portfolios obtained by applying (4) or (5) with total and partial ranking

Total ranking			Partial ranking		
id_{TR}	$w_i^{prop.TR}$	$w_i^{opt.TR}$	$w_i^{opt.PR}$	$w_i^{prop.PR}$	id_{PR}
4	0.462	0	0	0.299	142
63	0	0.581	0	0	34
130	0	0	0	0.225	109
49	0.003	0	0	0	49
109	0.216	0	1	0.344	159
159	0.204	0	0	0.009	65
85	0	0	0	0	63
54	0	0	0	0.123	147
60	0.004	0.419	0	0	4
34	0.112	0	0	0	28
Return	0.011	0.003	0.019	0.011	
St. dev.	1.185	0.77	0.655	0.913	
Sharpe	0.009	0.004	0.029	0.012	

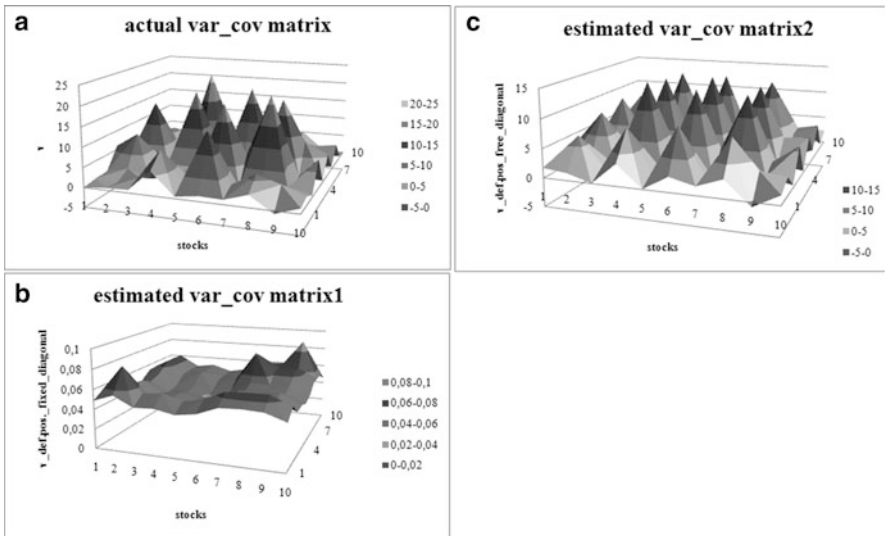


Fig. 2 Top left: actual var–cov matrix. Bottom left: estimated var–cov matrix fixed diagonal. Top right: estimated var–cov matrix free diagonal

There is empirical evidence that the portfolios obtained by removing the positiveness constrain of (5), applying (9), putting to zero the shorting weights and reportioning the others, give far better diversification.

(g) Approximations of the actual BEKK-sampling covariance matrix \hat{V} , see Fig. 2a, with the Higham algorithm are shown in Fig. 2b (fixed diagonal), Fig. 2c (free diagonal).

Although the fixed diagonal matrix has a better diagonal estimation given by the BEKK method part, the possibility to relax this constrain gives a better overall approximation.

- (h) The dual algorithm, starting from the unconstrained solution $\mathbf{0}$, converges in 12 iterations with all constraints active with the exception of the constraint 6 which is inactive but satisfied. The value of the objective function at the minimum is 0.0195.

References

1. Bollerslev, T., Engle R.F., Nelson D.B. "ARCH Model", Handbook of Econometrics, IV, pp. 2959–3038. Elsevier Science, Amsterdam (1994)
2. Bollerslev, T., Engle, R.F., Wooldridge, J.M.: A capital asset pricing model with time-varying covariance. *J. Polit. Econ.* **96**, 116–131 (1988)
3. Boyle, J.P., Dykstra, R.L.: A method for finding projections onto the intersection of convex sets in Hilbert spaces. In: *Advances in Order Restricted Inference. Lecture Notes in Statistics*, vol. 37, pp. 28–47. Springer, Berlin (1985)
4. Campbell, J.Y., Chan, Y.L., Viceira, L.M.: A multivariate model of strategic asset allocation. *J. Financ. Econ.* **67**, 41–80 (2003)
5. Datastream Global Equity Indices: User Guide, Issue 5. Thomson Reuters Ltd., Chennai (2008)
6. Engle, R.F., Kroner, K.F.: Multivariate simultaneous generalized ARCH. *Econ. Theor.* **11**, 122–150 (1995)
7. Evans-Archer: Diversification and the reduction of dispersion: an empirical analysis. *J. Financ.* **23**, 761–767 (1968)
8. Goldfarb, G., Idnani, A.: A numerical stable dual method for solving strictly convex quadratic programs. *Math. Prog.* **27**, 1–33 (1983)
9. Goodman, D.A., Peavy, J.W., III: Industry relative price-earnings ratios as indicators of investment returns. *Financ. Anal. J.* **39**(4), 60–66 (1983)
10. Higham, N.: Computing the nearest correlation matrix - a problem from finance. *IMA J. Numer. Anal.* **22**, 329–343 (2002)
11. Hlouskova, J., Schmidheiny, K., Wagner, M.: *Multistep Predictions from Multivariate ARMA-GARCH Models and their Value for Portfolio Management*. Universität Bern, Bern (2002)
12. Johansen, S.: *Likelihood-Based Inference in Cointegrated Vector Autoregressive models*. Oxford University Press, Oxford (1995)
13. Markowitz, H.: Portfolio selection. *J. Financ.* **7**, 77–91 (1952)