

Springer Proceedings in Mathematics & Statistics

Cira Perna · Monica Pratesi
Anne Ruiz-Gazen *Editors*

Studies in Theoretical and Applied Statistics

SIS 2016, Salerno, Italy, June 8–10



Springer Proceedings in Mathematics & Statistics

Volume 227

Springer Proceedings in Mathematics & Statistics

This book series features volumes composed of selected contributions from workshops and conferences in all areas of current research in mathematics and statistics, including operation research and optimization. In addition to an overall evaluation of the interest, scientific quality, and timeliness of each proposal at the hands of the publisher, individual contributions are all refereed to the high quality standards of leading journals in the field. Thus, this series provides the research community with well-edited, authoritative reports on developments in the most exciting areas of mathematical and statistical research today.

More information about this series at <http://www.springer.com/series/10533>

Cira Perna · Monica Pratesi
Anne Ruiz-Gazen
Editors

Studies in Theoretical and Applied Statistics

SIS 2016, Salerno, Italy, June 8–10



Editors

Cira Perna
Dipartimento di Scienze
Economiche e Statistiche
Università degli Studi di Salerno
Fisciano, Salerno
Italy

Anne Ruiz-Gazen
Toulouse School of Economics
University of Toulouse
Toulouse Cedex 6
France

Monica Pratesi
Dipartimento di Economia e
Management
Università degli Studi di Pisa
Pisa
Italy

ISSN 2194-1009 ISSN 2194-1017 (electronic)
Springer Proceedings in Mathematics & Statistics
ISBN 978-3-319-73905-2 ISBN 978-3-319-73906-9 (eBook)
<https://doi.org/10.1007/978-3-319-73906-9>

Library of Congress Control Number: 2018930101

Mathematics Subject Classification (2010): S11001, S17010, S12008, S17040

© Springer International Publishing AG, part of Springer Nature 2018

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. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by the registered company Springer International Publishing AG part of Springer Nature
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

This book contains a selection of the papers presented during the 48th Scientific Meeting of the Italian Statistical Society (SIS2016), held in Salerno on June 8–10, 2016.

This biennial conference is a traditional national and international meeting for connecting researchers in statistics, demography, and applied statistics in Italy. The conference aims at bringing together national and foreign researchers and practitioners to discuss recent developments in statistical methods for economics, social sciences, and all fields of application of statistics.

The Scientific Programme Committee provided a balanced and stimulating program that appealed to the diverse interests of the participants.

This book of selected papers is organized in chapters each related to a theme discussed in the meeting. In the editing process, we reordered the themes and collapsed some of them. The result still resembles the large variety of topics addressed in Salerno.

From the modern data sources and survey design issues to the study of the measures of sustainable development, the reader can find a large collection of research topics in Statistical Methods and in Applied Statistics and Demography.

In this respect, the papers included in this volume provide a comprehensive overview of the current Italian scientific researches in theoretical and applied statistics.

1. Open data and big data in public administration and official statistics
2. Survey sampling: theory and application
3. A recent debate in Statistics
4. Statistical algorithms
5. Ordinal and symbolic data
6. Statistical models and methods for network data
7. Forecasting time series
8. Spatial analysis
9. Issues on ecological and environmental statistics

10. Statistics and the education system
11. Economic and financial data analysis
12. Sustainable development: theory, measures, and applications

The Scientific Programme Committee, the Session Organizers, the local hosting University, and many volunteers have contributed substantially to the organization of the conference and to the referee process to obtain this book. We acknowledge their work and the support of our Society. Particularly, we wish to thank Marcella Niglio for her continuous support and assistance in the editing of this book.

Wishing you a productive and stimulating reading,

Salerno, Italy
Pisa, Italy
Toulouse, France
October 2017

Cira Perna
Monica Pratesi
Anne Ruiz-Gazen

Contents

Part I Advances in Survey Methods and New Sources in Public Statistics

**Robustness in Survey Sampling Using the Conditional Bias Approach
with R Implementation** 3

Cyril Favre-Martinoz, Anne Ruiz-Gazen, Jean Francois Beaumont
and David Haziza

**Methodological Perspectives for Surveying Rare and Clustered
Population: Towards a Sequentially Adaptive Approach** 15

Federico Andreis, Emanuela Furfaro and Fulvia Mecatti

**Age Management in Italian Companies. Findings
from Two INAPP Surveys** 25

Maria Laura Aversa, Paolo Emilio Cardone and Luisa D'Agostino

**Generating High Quality Administrative Data: New Technologies
in a National Statistical Reuse Perspective** 41

Manlio Calzaroni, Cristina Martelli and Antonio Samaritani

Exploring Solutions for Linking Big Data in Official Statistics 49

Tiziana Tuoto, Daniela Fusco and Loredana Di Consiglio

Part II Recent Debates in Statistics and Statistical Algorithms

An Algorithm for Finding Projections with Extreme Kurtosis 61

Cinzia Franceschini and Nicola Loperfido

**Maxima Units Search (MUS) Algorithm: Methodology
and Applications** 71

Leonardo Egidi, Roberta Pappadà, Francesco Pauli and Nicola Torelli

**DESPOTA: An Algorithm to Detect the Partition in the Extended
Hierarchy of a Dendrogram** 83

Davide Passaretti and Domenico Vistocco

| | |
|--|-----|
| The p-value Case, a Review of the Debate: Issues and Plausible Remedies | 95 |
| Francesco Pauli | |
| Part III Statistical Models and Methods for Network Data, Ordinal and Symbolic Data | |
| A Dynamic Discrete-Choice Model for Movement Flows | 107 |
| Johan Koskinen, Tim Müller and Thomas Grund | |
| On the Analysis of Time-Varying Affiliation Networks: The Case of Stage Co-productions | 119 |
| Giancarlo Ragozini, Marco Serino and Daniela D'Ambrosio | |
| Similarity and Dissimilarity Measures for Mixed Feature-Type Symbolic Data | 131 |
| Manabu Ichino and Kadri Umbleja | |
| Dimensionality Reduction Methods for Contingency Tables with Ordinal Variables | 145 |
| Luigi D'Ambra, Pietro Amenta and Antonello D'Ambra | |
| Part IV Forecasting Time Series | |
| Extended Realized GARCH Models | 159 |
| Richard Gerlach and Giuseppe Storti | |
| Updating CPI Weights Through Compositional VAR Forecasts: An Application to the Italian Index | 169 |
| Lisa Crosato and Biancamaria Zavanella | |
| Prediction Intervals for Heteroscedastic Series by Holt-Winters Methods | 179 |
| Paolo Chirico | |
| Part V Spatial Analysis and Issues on Ecological and Environmental Statistics | |
| Measuring Residential Segregation of Selected Foreign Groups with Aspatial and Spatial Evenness Indices. A Case Study | 189 |
| Federico Benassi, Frank Heins, Fabio Lipizzi and Evelina Paluzzi | |
| Space-Time FPCA Clustering of Multidimensional Curves | 201 |
| Giada Adelfio, Francesca Di Salvo and Marcello Chiodi | |
| The Power of Generalized Entropy for Biodiversity Assessment by Remote Sensing: An Open Source Approach | 211 |
| Duccio Rocchini, Luca Delucchi and Giovanni Bacaro | |

| | |
|--|-----|
| An Empirical Approach to Monitoring Ship CO₂ Emissions via Partial Least-Squares Regression | 219 |
| Antonio Lepore, Biagio Palumbo and Christian Capezza | |
| Part VI Statistics and the Education System | |
| Promoting Statistical Literacy to University Students: A New Approach Adopted by Istat | 231 |
| Alessandro Valentini, Monica Carbonara and Giulia De Candia | |
| From South to North? Mobility of Southern Italian Students at the Transition from the First to the Second Level University Degree | 239 |
| Marco Enea | |
| Monitoring School Performance Using Value-Added and Value-Table Models: Lessons from the UK | 251 |
| George Leckie and Harvey Goldstein | |
| Part VII Economic and Financial Data Analysis | |
| Indexing the Normalized Worthiness of Social Agents | 263 |
| Giulio D'Epifanio | |
| Financial Crises and Their Impacts: Data Gaps and Innovation in Statistical Production | 275 |
| Emanuele Baldacci | |
| European Welfare Systems in Official Statistics: National and Local Levels | 289 |
| Alessandra Coli and Barbara Pacini | |
| Financial Variables Analysis by Inequality Decomposition | 301 |
| Michele Costa | |
| Part VIII Sustainable Development: Theory, Measures and Applications | |
| A Novel Perspective in the Analysis of Sustainability, Inclusion and Smartness of Growth Through Europe 2020 Indicators | 311 |
| Elena Grimaccia and Tommaso Rondinella | |
| The Italian Population Behaviours Toward Environmental Sustainability: A Study from Istat Surveys | 325 |
| Isabella Mingo, Valentina Talucci and Paola Ungaro | |
| Estimating the at Risk of Poverty Rate Before and After Social Transfers at Provincial Level in Italy | 337 |
| Caterina Giusti and Stefano Marchetti | |

About the Editors

Cira Perna is currently Professor of Statistics and Head of the Department of Economics and Statistics, University of Salerno (Italy). Her research work mainly focuses on nonlinear time series, artificial neural network models, and resampling techniques. She has published a number of papers in national and international journals on these topics, and she has been a member of the scientific committees of several national and international conferences.

Monica Pratesi is Professor of Statistics, University of Pisa, and holds the Jean Monnet Chair “Small Area Methods for Monitoring of Poverty and Living Conditions in the EU” 2015–2017. She is the Director of the Tuscan Interuniversity Centre—Advanced Statistics for Equitable and Sustainable Development, entitled to Camilo Dagum. Her research interests include methods for survey sampling and analysis of survey data, small area estimation, and design-based population inference. She has published a number of papers in national and international journals on these topics and has been a member of the scientific committees of several national and international conferences.

Anne Ruiz-Gazen is Professor of Applied Mathematics, specializing in statistics, and a member of the Toulouse School of Economics—Research at University Toulouse 1 Capitole. Her areas of research include multivariate data analysis, survey sampling theory and, to a less extent, spatial econometrics and statistics. She has published more than fifty articles in refereed journals and books and has been a member of the scientific committees of several conferences.

Maxima Units Search (MUS) Algorithm: Methodology and Applications



Leonardo Egidi, Roberta Pappadà, Francesco Pauli
and Nicola Torelli

Abstract An algorithm for extracting identity submatrices of small rank and pivotal units from large and sparse matrices is proposed. The procedure has already been satisfactorily applied for solving the label switching problem in Bayesian mixture models. Here we introduce it on its own and explore possible applications in different contexts.

Keywords Identity matrix · Pivotal unit · Label switching

1 Introduction

Identifying and extracting identity matrices of small rank with given features from a larger, possibly sparse, matrix could appear just of theoretical interest. However, investigating the structure of a given sparse matrix is not only of theoretical appeal but can be useful for a wide variety of practical problems and for statistics.

This kind of matrix appears in clustering ensembles methods, which combine data partitions of the same dataset in order to obtain good data partitions even when the clusters are not compact and well separated. See, for instance, [1] where multiple partitions of the same data (an ensemble) are generated changing the number

L. Egidi (✉)

Dipartimento di Scienze Statistiche, Università degli Studi di Padova,
Via Cesare Battisti 241, 35121 Padova, Italy
e-mail: egidi@stat.unipd.it

R. Pappadà · F. Pauli · N. Torelli

Dipartimento di Scienze Economiche, Aziendali, Matematiche e Statistiche,
'Bruno de Finetti', Università degli Studi di Trieste, Via Tigor 22, 34124 Trieste, Italy
e-mail: rpappada@units.it

F. Pauli

e-mail: francesco.pauli@deams.units.it

N. Torelli

e-mail: nicola.torelli@deams.units.it

of clusters and using random cluster initializations within the K-means algorithm. Another situation where the global number of zeros of a matrix has a relevant role is in analysing the structure of a matrix of factor loadings; [4] introduces and formulates a statistical index in order to assess how good is the solution based on a factor analysis.

Matrices with a similar structure and for which the sparseness has to be taken into account appear in the so-called cost's optimization theory. [5] builds the well-known Hungarian method, which uses the zeros matrix elements for finding an optimal assignment for a given cost matrix; [6] presents a generalization of such algorithm and an application to a transportation problem.

In this paper we discuss the so-called Maxima Units Search algorithm (hereafter MUS). It has been introduced in [2] and used in the context of the label switching problem [3, 8]. In Bayesian estimation of finite mixture models label switching arises since the likelihood is invariant to permutations of the mixture components. The MUS procedure has proved to be useful in detecting some specific units—one for each mixture component—called pivots, from a large and sparse similarity matrix representing an estimate of the probability that pairs of units belong to the same group. The MUS algorithm is then more generally aimed at identifying for a given partition of the data those units that are not connected with a large number of units selected from the other groups.

A formal description of the MUS algorithm is provided and discussed. In fact, we argue that the proposed approach is of a broader interest and can be used for different purposes especially when the considered matrix presents a non-trivial number of zeros.

In Sect. 2 we introduce the notation, the algorithm and the main quantities of interest. A simulation study conducted for exploring the sensitivity of the algorithm to the choice of some parameters is presented in Sect. 3. Possible applications are illustrated in Sect. 4: in the first example we report the pivotal identification mentioned above, which represents the initial motivation for the procedure. Finally the method is employed to study a small dataset concerning tennis players' abilities. Section 5 concludes.

2 The Methodology

Let us consider a symmetric square matrix C of dimensions $N \times N$ containing a non-negligible number of zeros and suppose that each row's—or equivalently column's—index represents a statistical unit. Moreover, let us suppose that such N units either naturally belong to K different groups or have been preliminarily clustered into them, for instance via a suitable clustering technique.

For some practical purposes an example of which will be given in Sect. 4, we may be interested in detecting those units—one for each group—whose corresponding rows have more zeros than the other units. We preliminarily refer to these units as

the *maxima* units. More precisely, the underlying idea is to choose as maxima those units j_1, \dots, j_K such that the $K \times K$ submatrix of C , S_{j_1, \dots, j_K} with only the j_1, \dots, j_K rows and columns has few, possibly zero, non-zero elements off the diagonal (that is, the submatrix S_{j_1, \dots, j_K} is identical or nearly identical). Note that an identity submatrix of the given dimension may not exist. From a computational point of view, the issue is non-trivial and involves a global search row by row; as N , K and the number of zeros within C increase, the procedure becomes computationally demanding.

Before introducing mathematical details, let us denote with i_1, \dots, i_K a set of K maxima units and with S_{i_1, \dots, i_K} the submatrix of C containing only the rows and columns corresponding to the maxima. The main steps of the algorithm are summarized below.

- (i) For every group k , $k = 1, \dots, K$ find the *candidate maxima* units $j_k^1, \dots, j_k^{\bar{m}}$ within matrix C , i.e. the units in group k with the greater number of zeros corresponding to the units of the other $K - 1$ groups, where \bar{m} is a *precision parameter* fixed in advance. Let \mathcal{P}_k^h , $h = 1, \dots, \bar{m}$, $k = 1, \dots, K$ be the entire subset of units belonging to the remaining $K - 1$ groups which have a zero in j_k^h , that is

$$\mathcal{P}_k^h = \{j_l, l \neq k : C_{(j_k^h, j_l)} = 0\}, \quad h = 1, \dots, \bar{m}, k = 1, \dots, K$$

where $C_{(j_k^h, j_l)}$ is the element (j_k^h, j_l) of the matrix. We collect a total number of $\bar{m}K$ candidate maxima, \bar{m} for every group.

- (ii) For each of these $\bar{m}K$ units, count the number of distinct identity submatrices of C which contain them, constructed by taking a given candidate h and $K - 1$ elements of the corresponding set \mathcal{P}_k^h . Let us denote this quantity with

$$M_{j_k^h}^h = \#\{S_{j_1, \dots, j_{k-1}, j_k, j_{k+1}, \dots, j_K} \mid j_i \in \mathcal{P}_k^h, i = 1, \dots, k-1, k+1, \dots, K\}. \quad (1)$$

- (iii) For each group k , $k = 1, \dots, K$, select the unit which yields the greatest number of identity matrices of rank K . In mathematical terms

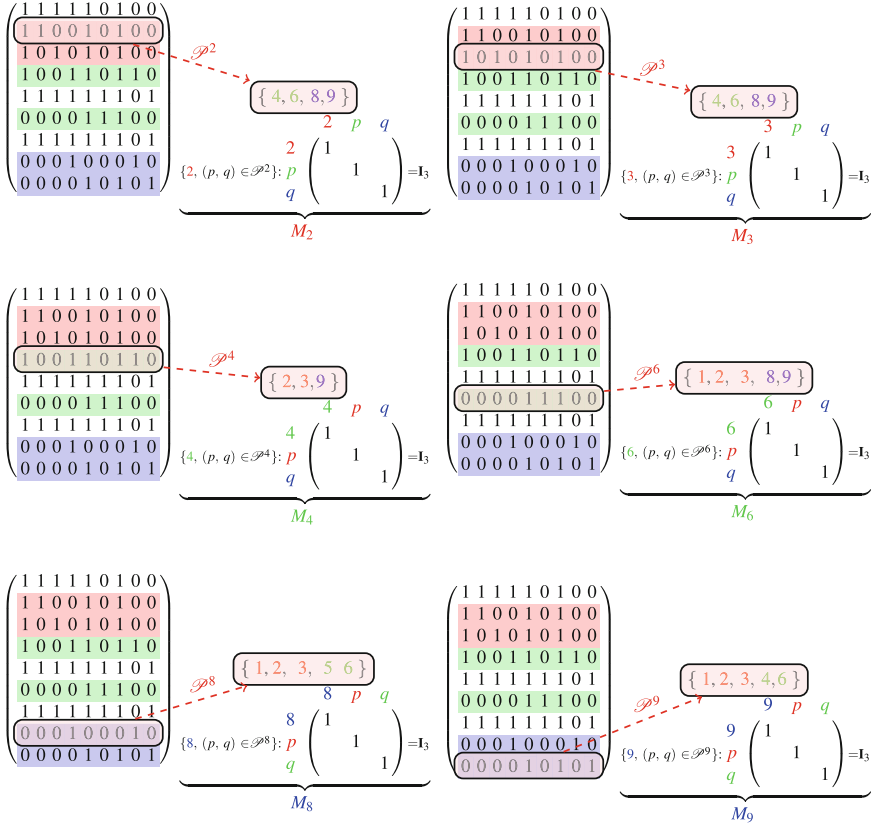
$$i_k = \operatorname{argmax}_{j_k^h \in \{j_k^1, \dots, j_k^{\bar{m}}\}} M_{j_k^h}^h, \quad h = 1, \dots, \bar{m}, k = 1, \dots, K. \quad (2)$$

The steps of the described algorithm are illustrated via a numerical example in Fig. 1. The choice of \bar{m} is crucial in terms of the algorithm performance. This parameter is a sort of benchmark for the size of the K subsets where the algorithm searches for the K maxima units: the greater is this value, the larger is the set of possible candidates involved in Eq. (2). Conversely, a bigger value enhances the possibility to build a larger set \mathcal{P}_k^h and obtain a more accurate result. In Sect. 3 we deal with this issue and we consider different choices for the precision parameter \bar{m} .

S. 1

$$C = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{matrix} & \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \end{matrix} \rightarrow \text{Candidates} = \{2, 3, 4, 6, 8, 9\}$$

S. 2



S. 3

$$\begin{aligned} \mathbf{i}_1 &= \textcircled{2} \text{ if } M_2 > M_3 \text{ or } \textcircled{3} \text{ if } M_3 > M_2 \\ \mathbf{i}_2 &= \textcircled{4} \text{ if } M_4 > M_6 \text{ or } \textcircled{6} \text{ if } M_6 > M_4 \\ \mathbf{i}_3 &= \textcircled{8} \text{ if } M_8 > M_9 \text{ or } \textcircled{9} \text{ if } M_9 > M_8 \end{aligned}$$

Fig. 1 Graphical scheme of the MUS algorithm for $K=3$ and precision parameter $\tilde{m}=2$. S. 1 Chooses the candidate maxima, the two units for each group with the greatest number of zeros. S. 2 Identifies for each candidate the subsets \mathcal{P} of units which belong to a different group (than the candidate) and have a zero in correspondence of it; then builds all the identity matrices of rank three which contain the candidates. S. 3 Detects the maxima as the three units—one for each group—that appear the greatest number of times in an identity matrix

3 Simulation Study

The task of this section is to investigate the performance of the MUS algorithm and its sensitivity to the choice of N and \bar{m} , for a fixed K , which is determined by some clustering technique or a given grouping of the units. To this aim, we simulate a symmetric $N \times N$ matrix C where the element (i, j) is drawn from a Bernoulli distribution with parameter p . As mentioned in Sect. 2, the i -th row's index, $i, i = 1, \dots, N$, is associated to a statistical unit of interest and each unit is here randomly assigned to group $k, k = 1, \dots, K$, with probability $1/K$. We consider three different values of p , i.e. $p = 0.8, 0.5, 0.2$.

Tables 1, 2 and 3 display the maxima units and the corresponding CPU times in seconds (in brackets) according to the considered scenarios. As expected, the procedure is sensitive to the choices of input weights, both in terms of units selection and computational times.

The first issue one may immediately notice is that, regardless of the weights used for generating data, the computational burden rises dramatically when $K > 3$. Especially when $N = 1000$, the CPU time is huge if compared to the time spent in the same framework—same \bar{m} and weights—for $K = 2$ or $K = 3$. As the probability p decreases (from 0.8 to 0.2) the number of zeros becomes larger and, consequently, the CPU time required keeps growing regardless of the values of N and \bar{m} . A second remark is that, by fixing N and K , the choice of the precision parameter \bar{m} does not seem to affect significantly the performance of the procedure: as \bar{m} increases, there is limited variation in the units detection and the difference in the required time between $\bar{m} = 1$ and $\bar{m} = 20$ remains relatively small, as can be seen from Tables 1, 2 and 3. This is suggesting that even the choice of a small precision parameter—e.g. $\bar{m} = 5$ —may be accurate enough for detecting the maxima.

4 Applications

4.1 Identification of Pivotal Units

As broadly explained in [2], the identification of some pivotal units in a Bayesian mixture model with a fixed number of groups may be helpful when dealing with the label switching problem [3, 8].

Let N be the number of observations generated from the mixture model. Consider, for instance, the probability of two units being in the same group. Such quantity may be estimated from the MCMC sample and denoted as \hat{c}_{ij} . For details, see [2]. The $N \times N$ matrix C with elements \hat{c}_{ij} can be seen as a similarity matrix between units. Now, such matrix can be considered as input for some suitable clustering techniques, in order to obtain a partition of the N observations into K groups. From such partition, we may be interested in identifying exactly K pivotal units—*pivots*—which are (pair-wise) separated with (posterior) probability one (that is, the posterior probability of

Table 1 MUS algorithm's maxima and computational times (in brackets) according to $K = 2, 3, 4$, $N = 100, 500, 1000$, $\bar{m} = 1, 5, 10, 20$. Bernoulli data 0, 1 generated with weights $p = 0.8$

| | | $K = 2$ | $K = 3$ | $K = 4$ |
|------------|----------------|-----------------|----------------------|-----------------------------|
| $N = 100$ | $\bar{m} = 1$ | 43, 14 (< 0.1) | -- -- (< 0.1) | 18, 34, 41, 88 (< 0.1) |
| | $\bar{m} = 5$ | 16, 14 (< 0.1) | 10, 49, 96 (< 0.1) | 37, 78, 17, 69 (0.16) |
| | $\bar{m} = 10$ | 16, 14 (< 0.1) | 10, 49, 96 (< 0.1) | 37, 78, 65, 69 (0.25) |
| | $\bar{m} = 20$ | 16, 14 (< 0.1) | 10, 49, 96 (< 0.1) | 37, 78, 65, 69 (0.43) |
| $N = 500$ | $\bar{m} = 1$ | -- -- (0.56) | -- -- (0.63) | 27, 44, 59, 263, (2.3) |
| | $\bar{m} = 5$ | 183, 125 (0.66) | 346, 373, 500 (0.68) | 394, 44, 59, 263, (9.7) |
| | $\bar{m} = 10$ | 183, 125 (0.64) | 399, 373, 500 (0.94) | 394, 44, 59, 263, (17.6) |
| | $\bar{m} = 20$ | 183, 125 (0.72) | 399, 373, 500 (1.22) | 394, 44, 59, 263, (32.71) |
| $N = 1000$ | $\bar{m} = 1$ | -- -- (2.32) | -- -- (2.40) | 350, 825, 916, 204 (10.9) |
| | $\bar{m} = 5$ | 654, 94 (2.49) | 909, 499, 868 (3.02) | 381, 849, 684, 204 (44.0) |
| | $\bar{m} = 10$ | 654, 94 (2.62) | 909, 499, 868 (3.52) | 381, 849, 684, 488 (81.7) |
| | $\bar{m} = 20$ | 654, 96 (2.99) | 909, 382, 868 (4.62) | 381, 849, 748, 488 (152.82) |

Table 2 MUS algorithm's maxima and computational times (in brackets) according to $K = 2, 3, 4$, $N = 100, 500, 1000$, $\bar{m} = 1, 5, 10, 20$, Bernoulli data 0, 1 generated with weights $p = 0.5$

| | | $K = 2$ | $K = 3$ | $K = 4$ |
|------------|----------------|--------------------------|--------------------------------|---------------------------------------|
| $N = 100$ | $\bar{m} = 1$ | 48, 86 (< 0.1) | -- -- (< 0.1) | 32, 62, 38, 55 (0.44) |
| | $\bar{m} = 5$ | 48, 61 (< 0.1) | 15, 33, 5, (0.12) | 50, 62, 89, 55 (1.99) |
| | $\bar{m} = 10$ | 48, 61 (< 0.1) | 15, 62, 5 (0.14) | 50, 62, 90, 55 (3.36) |
| | $\bar{m} = 20$ | 48, 61 (0.10) | 15, 62, 5 (0.13) | 50, 62, 90, 55 (1328.73) |
| $N = 500$ | $\bar{m} = 1$ | -- -- (1.61) | -- -- -- (1.67) | 10, 11, 90, 488 (56.1) |
| | $\bar{m} = 5$ | 294, 242 (1.40) | 203, 213, 272 (2.31) | 273, 242, 292, 383 (159.8) |
| | $\bar{m} = 10$ | 294, 242 (1.64) | 203, 213, 272 (3.31) | 273, 232, 482, 383 (311.28) |
| | $\bar{m} = 20$ | 66, 242 (1.78) | 203, 213, 272 (5.49) | 273, 232, 29, 383 (582.38) |
| $N = 1000$ | $\bar{m} = 1$ | -- -- (6.64) | -- -- -- (7.16) | 123, 964, 813, 238 (246.28) |
| | $\bar{m} = 5$ | 94, 405 (6.81) | 67, 995, 688, (9.78) | 267, 964, 813, 241 (1208.27) |
| | $\bar{m} = 10$ | 94, 405 (7.24) | 67, 995, 688, (12.61) | 267, 964, 813, 241 (2326.64) |
| | $\bar{m} = 20$ | 398, 405 (8.23) | 67, 995, 688, (9.58) | 267, 964, 813, 241 (4548.47) |

Table 3 MUS algorithm's maxima and computational times (in brackets) according to $K = 2, 3, 4$, $N = 100, 500, 1000$, $\bar{m} = 1, 5, 10, 20$. Bernoulli data 0, 1 generated with weights $p = 0.2$

| | | $K = 2$ | $K = 3$ | $K = 4$ |
|------------|----------------|-------------------|-----------------------|-------------------------------|
| $N = 100$ | $\bar{m} = 1$ | 42, 32 (< 0.1) | 58, 40, 63 (0.10) | 24, 68, 34, 89 (3.75) |
| | $\bar{m} = 5$ | 42, 32 (< 0.1) | 81, 54, 63 (0.21) | 24, 68, 48, 58 (8.85) |
| | $\bar{m} = 10$ | 42, 86 (0.14) | 87, 54, 63 (0.25) | 24, 68, 48, 58 (16.8) |
| | $\bar{m} = 20$ | 42, 86 (0.11) | 87, 54, 63 (0.43) | 24, 68, 48, 58 (1985.01) |
| $N = 500$ | $\bar{m} = 1$ | -- (2.40) | -- -- (2.74) | 371, 28, 122, 60 (189.94) |
| | $\bar{m} = 5$ | 326, 288 (2.39) | 290, 393, 316 (4.51) | 370, 38, 202, 404 (949.4) |
| | $\bar{m} = 10$ | 326, 288 (2.65) | 290, 393, 316 (7.20) | 370, 413, 202, 196 (1882.93) |
| | $\bar{m} = 20$ | 284, 288 (3.06) | 375, 395, 316 (10.66) | 370, 38, 202, 404 (3685.61) |
| $N = 1000$ | $\bar{m} = 1$ | 555, 892, (11.30) | -- -- -- (12.25) | 427, 452, 218, 631 (1608.11) |
| | $\bar{m} = 5$ | 434, 892 (11.28) | 222, 921, 275 (19.42) | 427, 493, 218, 839 (8098.54) |
| | $\bar{m} = 10$ | 434, 892 (11.27) | 387, 921, 255 (26.72) | 427, 493, 218, 839 (16629.86) |
| | $\bar{m} = 20$ | 434, 892 (12.47) | 387, 921, 255 (45.08) | 427, 493, 218, 839 (32230.8) |

any two of them being in the same group is zero). In fact, as discussed in [2], the identification of such units allows us to provide a valid solution to the occurrence of label switches.

Following the procedure described in Sect. 2, one can find units i_1, \dots, i_K such that the submatrix S of C , with only the rows and columns corresponding to such units, is the identity matrix. It is still worth noticing that the availability of K perfectly separated units is crucial to the procedure, and it can not always be guaranteed.

Practically, our interest is in finding units which should be ‘as far as possible’ one from each other according to a well defined distance measure. The more separated they are, the better they represent the group they belong to.

Figure 2 shows the pivotal identification of $K = 4$ units for a sample of $N = 1000$ bivariate data generated according to a nested Gaussian mixture of mixtures with K groups and fixed means. Pivots (red) have been detected through the MUS algorithm and through another alternative method, which aims at searching the most distant units among the members that are farthest apart. We may graphically notice that separation is made more efficient by the MUS algorithm, for which the red points appear quite distant from each other. Moreover, in this specific example the pivotal search is made difficult due to the overlapping of the K groups.

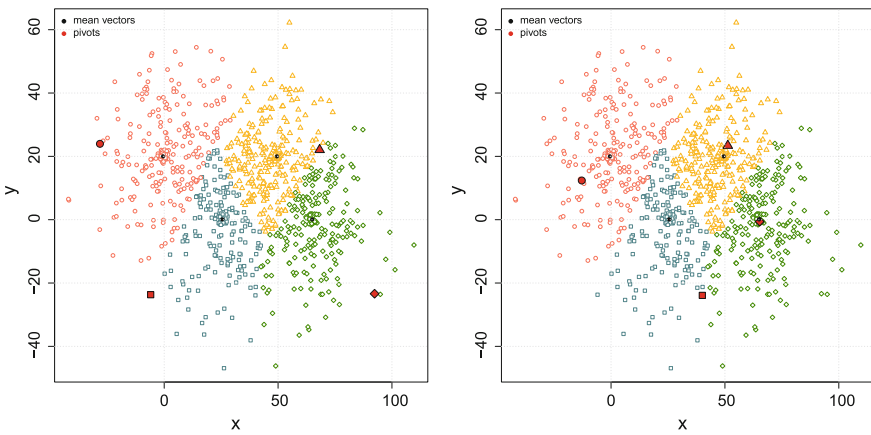


Fig. 2 Simulated bivariate sample of size $N = 1000$ from a nested Gaussian mixture of mixtures with $K = 4$ and input means (in black) $\mu_1 = (25, 0)$, $\mu_2 = (60, 0)$, $\mu_3 = (0, 20)$, $\mu_4 = (50, 20)$. Groups have been detected through an agglomerative clustering technique. Pivots—i.e. maxima—found by MUS algorithm (Left) with $\bar{m} = 5$ are shown in red and seem well separated in the bi-dimensional space. Pivots found by method $\min_i(\min_{j \notin \mathcal{X}_k} C_{ij})$ are less distant each other (Right)

4.2 Tennis Singular Features

As a simple example we apply our algorithm to a case study regarding tennis players. We collect $N = 8$ game’s features (hereafter GF) for $T = 25$ players from the Wimbledon Tournament 2016,¹ and we assign the following values

$$\begin{cases} GF_{i,t} = 1, & \text{if player } t \text{ has } GF_i, \\ GF_{i,t} = 0, & \text{if player } t \text{ has not the } GF_i. \end{cases}$$

Game’s features belong to $K = 2$ groups, which somehow refer to the attack and the defence skills for each player. We denote with the label 1 the first group and with 2 the second group: ‘First Serve Receiving Points’ (2), ‘Second Serve Receiving Points’ (2), ‘Break Points Won’ (1), ‘Serve Speed’ (1), ‘Aces’(1), ‘First Serve Points’ (1), ‘Second Serve Points’ (1), ‘Break Points Conversion’ (2).

We decide to assign a specific game feature to a given player if this player belongs to the first five positions of that game’s feature rank reported by the Wimbledon’s website. Hence, let us enumerate the above-mentioned features from one to eight. Our 0–1 dataset has as many records as players. The problem setting is summarized below.

| | GF | | | | | | | |
|---------|----|---|---|---|---|---|---|---|
| Player | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| Federer | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Murray | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |

Note that Federer is assigned game’s feature one (‘First Serve Receiving Points’) and two (‘Second Serve Receiving Points’) only, Murray is assigned game’s feature one, two and three (‘Break Points Won’) and so on. We define the $N \times N$ symmetric matrix, C , in which the generic element $C_{(i,j)}$ is the number of players that have both features i and j .

$$C = \begin{pmatrix} 1 & 5 & 3 & 1 & 3 & 0 & 0 & 0 \\ 5 & 1 & 2 & 1 & 2 & 0 & 0 & 0 \\ 3 & 2 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 2 & 0 & 0 & 0 \\ 3 & 2 & 1 & 2 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

We would like to extract the two ‘most distant’ game’s features for the two groups, i.e. the two features in correspondence of which the matrix C is zero more often, and

¹http://www.wimbledon.com/en_GB/scores/extrastats/index.html.

for which S_{i_1, i_2} is an identity matrix. We can notice that rows 6 and 7 of C are full of zeros: this means, for instance, that according to our short dataset $C_{(6,1)} = 0$, i.e. the ‘First Serve Points’ ($k = 1$) and ‘First Serve Receiving Points’ ($k = 2$) do not coexist to any player. Are they the most distant features between the two groups? To answer this question, we run the MUS algorithm by fixing $\bar{m} = 3$ and we find the candidate maxima $j_1^1 = 6$, $j_1^2 = 7$, $j_2^1 = 8$, $j_2^2 = 1$ and maxima $i_1 = 6$, $i_2 = 8$. Hence we conclude that ‘First Serve Points’ (six) and ‘Break Points Conversion’ (eight) are quite unlikely to belong to the same player.

5 Conclusions

A procedure for detecting small identity submatrices from a $N \times N$ matrix has been proposed. It has been initially considered for application to the pivotal approach in label switching problem in the analysis of Bayesian mixture models. The proposed method is discussed in detail and employed for different practical problems.

Its efficiency and its sensitivity to parameter choices is investigated through a simulation study, which shows that for a small number of groups the procedure is quite fast. Moreover, even for small values of the precision parameter \bar{m} the procedure appears quite stable in terms of units indexes, suggesting that a higher value of \bar{m} is often not required. This is also confirmed by the results in Sect. 4.

Further issues for future research are related to the optimization of the proposed algorithm and the definition of suitable indicators for detecting both diagnostic problems inherent to the procedure and goodness of units choice.

References

1. Ana, L.N.F., Jain, A.K.: Data clustering using evidence accumulation. In: 2002 Proceedings 16th International Conference on Pattern Recognition, vol. 4, pp. 276–280 (2002)
2. Egidi, L., Pappadà, R., Pauli, F., Torelli, N.: Relabelling in Bayesian mixture models by pivotal units. [arXiv:1501.05478v2](https://arxiv.org/abs/1501.05478v2) (2015)
3. Jasra, A., Holmes, C.C., Stephens, D.A.: Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling. *Stat. Sci.* 50–67 (2005)
4. Kaiser, Henry F.: An index of factorial simplicity. *Psychometrika* **39**(1), 31–36 (1974)
5. Kuhn, H.W.: The Hungarian method for the assignment problem. *Nav. Res. Logist. Q.* **2**(1–2), 83–97 (1955)
6. Munkres, J.: Algorithms for the assignment and transportation problems. *J. Soc. Ind. Appl. Math.* **5**(1), 32–38 (1957)
7. Puolamäki, K., Kaski, S.: Bayesian solutions to the label switching problem. In: *Advances in Intelligent Data Analysis VIII*, pp. 381–392. Springer (2009)
8. Stephens, M.: Dealing with label switching in mixture models. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **62**(4), 795–809 (2000)