

DIPARTIMENTO DI ECONOMIA E GIURISPRUDENZA
UNIVERSITÀ DI CASSINO E DEL LAZIO MERIDIONALE



CLADAG 2019

11-13 SEPTEMBER 2019
CASSINO

```
def business_model()  
  arr = [ ]  
  items = a, b, c  
  items >> arr  
  return arr  
end
```



Book of Short Papers

Giovanni C. Porzio
Francesca Greselin
Simona Balzano
Editors

12-TH SCIENTIFIC MEETING
CLASSIFICATION AND DATA ANALYSIS

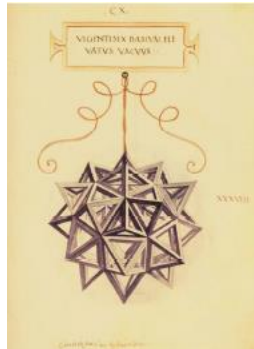


© CC – Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)
<https://creativecommons.org/licenses/by-nc/4.0/>

2019

Università di Cassino e del Lazio Meridionale
Centro Editoriale di Ateneo
Palazzo degli Studi Località Folcara, Cassino (FR), Italia

ISBN 978-88-8317-108-6



CLADAG 2019
Book of Short Papers

Giovanni C. Porzio
Francesca Greselin
Simona Balzano
Editors

2019

Contents

Keynotes lectures

Unifying data units and models in (co-)clustering <i>Christophe Biernacki</i>	3
Statistics with a human face <i>Adrian Bowman</i>	4
Bayesian model-based clustering with flexible and sparse priors <i>Bettina Grün</i>	5
Grinding massive information into feasible statistics: current challenges and opportunities for data scientists <i>Francesco Mola</i>	6
Statistical challenges in the analysis of complex responses in biomedicine <i>Sylvia Richardson</i>	7

Invited and contributed sessions

Model-based clustering of time series data: a flexible approach using nonparametric state-switching quantile regression models <i>Timo Adam, Roland Langrock, Thomas Kneib</i>	8
Some issues in generalized linear modeling <i>Alan Agresti</i>	12
Assessing social interest in burnout using functional data analysis through google trends <i>Ana M. Aguilera, Francesca Fortuna, Manuel Escabias</i>	16
Measuring equitable and sustainable well-being in Italian regions: a non- aggregative approach <i>Leonardo Salvatore Alaimo, Filomena Maggino</i>	20
Bootstrap inference for missing data reconstruction <i>Giuseppina Albano, Michele La Rocca, Maria Lucia Parrella, Cira Perna</i>	22
Archetypal contour shapes <i>Aleix Alcacer, Irene Epifanio, M. Victoria Ibáñez, Amelia Simó</i>	26

Random projections of variables and units <i>Laura Anderlucci, Roberta Falcone, Angela Montanari</i>	30
Sparse linear regression via random projections ensembles <i>Laura Anderlucci, Matteo Farnè, Giuliano Galimberti, Angela Montanari</i>	34
High-dimensional model-based clustering via random projections <i>Laura Anderlucci, Francesca Fortunato, Angela Montanari</i>	38
Multivariate outlier detection in high reliability standards fields using ICS <i>Aurore Archimbaud, Klaus Nordhausen, Anne Ruiz-Gazen</i>	42
Evaluating the school effect: adjusting for pre-test or using gain scores? <i>Bruno Arpino, Silvia Bacci, Leonardo Grilli, Raffaele Guetto, Carla Rampichini</i>	45
ACE, AVAS and robust data transformations <i>Anthony Atkinson</i>	49
Mixtures of multivariate leptokurtic Normal distributions <i>Luca Bagnato, Antonio Punzo, Maria Grazia Zoia</i>	53
Detecting and interpreting the consensus ranking based on the weighted Kemeny distance <i>Alessio Baldassarre, Claudio Conversano, Antonio D'Ambrosio</i>	57
Predictive principal components analysis <i>Simona Balzano, Maja Bozic, Laura Marcis, Renato Salvatore</i>	61
Flexible model-based trees for count data <i>Federico Banchelli</i>	63
Euclidean distance as a measure of conformity to Benford's law in digital analysis for fraud detection <i>Mateusz Baryła, Józef Pociecha</i>	67
The evolution of the purchase behavior of sparkling wines in the Italian market <i>Francesca Bassi, Fulvia Pennoni, Luca Rossetto</i>	71
Modern likelihood-frequentist inference at work <i>Ruggero Bellio, Donald A. Pierce</i>	75
Ontology-based classification of multilingual corpuses of documents <i>Sergey Belov, Salvatore Ingrassia, Zoran Kalinić, Paweł Lula</i>	79
Modeling heterogeneity in clustered data using recursive partitioning <i>Moritz Berger, Gerhard Tutz</i>	83

Mixtures of experts with flexible concomitant covariate effects: a bayesian solution <i>Marco Berrettini, Giuliano Galimberti, Thomas Brendan Murphy, Saverio Ranciati</i>	87
Sampling properties of an ordinal measure of interrater absolute agreement <i>Giuseppe Bove, Pier Luigi Conti, Daniela Marella</i>	91
Tensor analysis can give better insight <i>Rasmus Bro</i>	95
A boxplot for spherical data <i>Davide Buttarazzi, Giuseppe Pandolfo, Giovanni C. Porzio, Christophe Ley</i>	97
Machine learning models for forecasting stock trends <i>Giacomo Camba, Claudio Conversano</i>	99
Tree modeling ordinal responses: CUBREMOT and its applications <i>Carmela Cappelli, Rosaria Simone, Francesca Di Iorio</i>	103
Supervised learning in presence of outliers, label noise and unobserved classes <i>Andrea Cappelozzo, Francesca Greselin, Thomas Brendan Murphy</i>	104
Asymptotics for bandwidth selection in nonparametric clustering <i>Alessandro Casa, José E. Chacón, Giovanna Menardi</i>	108
Foreign immigration and pull factors in Italy: a spatial approach <i>Oliviero Casacchia, Luisa Natale, Francesco Giovanni Truglia</i>	112
Dimensionality reduction via hierarchical factorial structure <i>Carlo Cavicchia, Maurizio Vichi, Giorgia Zaccaria</i>	116
Likelihood-type methods for comparing clustering solutions <i>Luca Coraggio, Pietro Coretto</i>	120
Labour market analysis through transformations and robust multilevel models <i>Aldo Corbellini, Marco Magnani, Gianluca Morelli</i>	124
Modelling consumers' qualitative perceptions of inflation <i>Marcella Corduas, Rosaria Simone, Domenico Piccolo</i>	128
Noise resistant clustering of high-dimensional gene expression data <i>Pietro Coretto, Angela Serra, Roberto Tagliaferri</i>	132
Classify X-ray images using convolutional neural networks <i>Federica Crobu, Agostino Di Ciaccio</i>	136

A compositional analysis approach assessing the spatial distribution of trees in Guadalajara, Mexico <i>Marco Antonio Cruz, Maribel Ortego, Elisabet Roca</i>	140
Joining factorial methods and blockmodeling for the analysis of affiliation networks <i>Daniela D'Ambrosio, Marco Serino, Giancarlo Ragozini</i>	142
A latent space model for clustering in multiplex data <i>Silvia D'Angelo, Michael Fop</i>	146
Post processing of two dimensional road profiles: variogram scheme application and sectioning procedure <i>Mauro D'Apuzzo, Rose-Line Spacagna, Azzurra Evangelisti, Daniela Santilli, Vittorio Nicolosi</i>	150
A new approach to preference mapping through quantile regression <i>Cristina Davino, Tormod Naes, Rosaria Romano, Domenico Vistocco</i>	154
On the robustness of the cosine distribution depth classifier <i>Houyem Demni, Amor Messaoud, Giovanni C. Porzio</i>	158
Network effect on individual scientific performance: a longitudinal study on an Italian scientific community <i>Domenico De Stefano, Giuseppe Giordano, Susanna Zaccarin</i>	162
Penalized vs constrained maximum likelihood approaches for clusterwise linear regression modelling <i>Roberto Di Mari, Stefano Antonio Gattone, Roberto Rocci</i>	166
Local fitting of angular variables observed with error <i>Marco Di Marzio, Stefania Fensore, Agnese Panzera, Charles C. Taylor</i>	170
Quantile composite-based path modeling to estimate the conditional quantiles of health indicators <i>Pasquale Dolce, Cristina Davino, Stefania Taralli, Domenico Vistocco</i>	174
AUC-based gradient boosting for imbalanced classification <i>Martina Dossi, Giovanna Menardi</i>	178
How to measure material deprivation? A latent Markov model based approach <i>Francesco Dotto</i>	182
Decomposition of the interval based composite indicators by means of biclustering <i>Carlo Drago</i>	186
Consensus clustering via pivotal methods <i>Leonardo Egidi, Roberta Pappadà, Francesco Pauli, Nicola Torelli</i>	190

Robust model-based clustering with mild and gross outliers <i>Alessio Farcomeni, Antonio Punzo</i>	194
Gaussian processes for curve prediction and classification <i>Sara Fontanella, Lara Fontanella, Rosalba Ignaccolo, Luigi Ippoliti, Pasquale Valentini</i>	198
A new proposal for building immigrant integration composite indicator <i>Mario Fordellone, Venera Tomaselli, Maurizio Vichi</i>	199
Biodiversity spatial clustering <i>Francesca Fortuna, Fabrizio Maturo, Tonio Di Battista</i>	203
Skewed distributions or transformations? Incorporating skewness in a cluster analysis <i>Michael Gallagher, Paul McNicholas, Volodymyr Melnykov, Xuwen Zhu</i>	207
Robust parsimonious clustering models <i>Luis Angel Garcia-Escudero, Agustin Mayo-Isacar, Marco Riani</i>	208
Projection-based uniformity tests for directional data <i>Eduardo García-Portugués, Paula Navarro-Esteban, Juan Antonio Cuesta-Albertos</i>	212
Graph-based clustering of visitors' trajectories at exhibitions <i>Martina Gentilin, Pietro Lovato, Gloria Menegaz, Marco Cristani, Marco Minozzo</i>	214
Symmetry in graph clustering <i>Andreas Geyer-Schulz, Fabian Ball</i>	218
Bayesian networks for the analysis of entrepreneurial microcredit: evidence from Italy <i>Lorenzo Giammei, Paola Vicard</i>	222
The PARAFAC model in the maximum likelihood approach <i>Paolo Giordani, Roberto Rocci, Giuseppe Bove</i>	226
Structure discovering in nonparametric regression by the GRID procedure <i>Francesco Giordano, Soumendra Nath Lahiri, Maria Lucia Parrella</i>	230
A microblog auxiliary part-of-speech tagger based on bayesian networks <i>Silvia Golia, Paola Zola</i>	234
Recent advances in model-based clustering of high dimensional data <i>Isobel Claire Gormley</i>	238
Tree embedded linear mixed models <i>Anna Gottard, Leonardo Grilli, Carla Rampichini, Giulia Vannucci</i>	239

Weighted likelihood estimation of mixtures <i>Luca Greco, Claudio Agostinelli</i>	243
A canonical representation for multiblock methods <i>Mohamed Hanafi</i>	247
An adequacy approach to estimating the number of clusters <i>Christian Hennig</i>	251
Classification with weighted compositions <i>Karel Hron, Julie Rendlova, Peter Filzmoser</i>	255
MacroPCA: an all-in-one PCA method allowing for missing values as well as cellwise and rowwise outliers <i>Mia Hubert, Peter J. Rousseeuw, Wannes Van den Bossche</i>	256
Marginal effects for comparing groups in regression models for ordinal outcome when uncertainty is present <i>Maria Iannario, Claudia Tarantola</i>	258
A multi-criteria approach in a financial portfolio selection framework <i>Carmela Iorio, Giuseppe Pandolfo, Roberta Siciliano</i>	262
Clustering of trajectories using adaptive distances and warping <i>Antonio Irpino, Antonio Balzanella</i>	266
Sampling and learning Mallows and generalized Mallows models under the Cayley distance: short paper <i>Ekhine Irurozki, Borja Calvo, Jose A. Lozano</i>	270
The gender parity index for the academic students progress <i>Aglaia Kalamatianou, Adele H. Marshall, Mariangela Zenga</i>	274
Some asymptotic properties of model selection criteria in the latent block model <i>Christine Keribin</i>	278
Invariant concept classes for transcriptome classification <i>Hans Kestler, Robin Szekely, Attila Klimmek, Ludwig Lausser</i>	282
Clustering of ties defined as symbolic data <i>Luka Kronegger</i>	283
Application of data mining in the housing affordability analysis <i>Viera Labudová, Eubica Sipková</i>	284
Cylindrical hidden Markov fields <i>Francesco Lagona</i>	288

Comparing tree kernels performances in argumentative evidence classification <i>Davide Liga</i>	292
Recent advancement in neural network analysis of biomedical big data <i>Pietro Liò, Giovanna Maria Dimitri, Chiara Sopegno</i>	296
Bias reduction for estimating functions and pseudolikelihoods <i>Nicola Lunardon</i>	297
Large scale social and multilayer networks <i>Matteo Magnani</i>	301
Uncertainty in statistical matching by BNs <i>Daniela Marella, Paola Vicard, Vincenzina Vitale</i>	305
Evaluating the recruiters' gender bias in graduate competencies <i>Paolo Mariani, Andrea Marletta</i>	309
Dynamic clustering of network data: a hybrid maximum likelihood approach <i>Maria Francesca Marino, Silvia Pandolfi</i>	313
Stability of joint dimension reduction and clustering <i>Angelos Markos, Michel Van de Velden, Alfonso Iodice D'Enza</i>	317
Hidden Markov models for clustering functional data <i>Andrea Martino, Giuseppina Guatteri, Anna Maria Paganoni</i>	321
Composite likelihood inference for simultaneous clustering and dimensionality reduction of mixed-type longitudinal data <i>Antonello Maruotti, Monia Ranalli, Roberto Rocci</i>	325
Bivariate semi-parametric mixed-effects models for classifying the effects of Italian classes on multiple student achievements <i>Chiara Masci, Francesca Ieva, Tommaso Agasisti, Anna Maria Paganoni</i>	329
Multivariate change-point analysis for climate time series <i>Gianluca Mastrantonio, Giovanna Jona Lasinio, Alessio Pollice, Giulia Capotorti, Lorenzo Teodonio, Carlo Blasi</i>	333
A dynamic stochastic block model for longitudinal networks <i>Catherine Matias, Tabea Rebafka, Fanny Villers</i>	337
Unsupervised fuzzy classification for detecting similar functional objects <i>Fabrizio Mauro, Francesca Fortuna, Tonio Di Battista</i>	339
Mixture modelling with skew-symmetric component distributions <i>Geoffrey McLachlan</i>	343

New developments in applications of pairwise overlap <i>Volodymyr Melnykov, Yana Melnykov, Domenico Perrotta, Marco Riani, Francesca Torti, Yang Wang</i>	344
Modelling unobserved heterogeneity of ranking data with the bayesian mixture of extended Plackett-Luce models <i>Cristina Mollica, Luca Tardella</i>	346
Issues in nonlinear time series modeling of European import volumes <i>Gianluca Morelli, Francesca Torti</i>	350
Gaussian parsimonious clustering models with covariates and a noise component <i>Keefe Murphy, Thomas Brendan Murphy</i>	352
Illumination in depth analysis <i>Stanislav Nagy, Jiří Dvořák</i>	353
Copula-based non-metric unfolding on augmented data matrix <i>Marta Nai Ruscone, Antonio D'Ambrosio</i>	357
A statistical model for software releases complexity prediction <i>Marco Ortu, Giuseppe Destefanis, Roberto Tonelli</i>	361
Comparison of serious diseases mortality in regions of V4 <i>Viera Pacáková, Lucie Kopecká</i>	365
Price and product design strategies for manufacturers of electric vehicle batteries: inferences from latent class analysis <i>Friederike Paetz</i>	369
A Mahalanobis-like distance for cylindrical data <i>Lucio Palazzo, Giovanni C. Porzio, Giuseppe Pandolfo</i>	373
Archetypes, prototypes and other types <i>Francesco Palumbo, Giancarlo Ragozini, Domenico Vistocco</i>	377
Generalizing the skew-t model using copulas <i>Antonio Parisi, Brunero Liseo</i>	381
Contamination and manipulation of trade data: the two faces of customs fraud <i>Domenico Perrotta, Andrea Cerasa, Lucio Barabesi, Mario Menegatti, Andrea Cerioli</i>	385
Bayesian clustering using non-negative matrix factorization <i>Michael Porter, Ketong Wang</i>	389

Exploring gender gap in international mobility flows through a network analysis approach <i>Ilaria Primerano, Marialuisa Restaino</i>	393
Clustering two-mode binary network data with overlapping mixture model and covariates information <i>Saverio Ranciati, Veronica Vinciotti, Ernst C. Wit, Giuliano Galimberti</i>	395
A stochastic blockmodel for network interaction lengths over continuous time <i>Riccardo Rastelli, Michael Fop</i>	399
Computationally efficient inference for latent position network models <i>Riccardo Rastelli, Florian Maire, Nial Friel</i>	403
Clustering of complex data stream based on barycentric coordinates <i>Parisa Rastin, Basarab Matei, Guénaél Cabanes</i>	407
An INDSCAL based mixture model to cluster mixed-type of data <i>Roberto Rocci, Monia Ranalli</i>	411
Topological stochastic neighbor embedding <i>Nicoleta Rogovschi, Nistor Grozavu, Basarab Matei, Younès Bennani, Seiichi Ozawa</i>	415
Functional data analysis for spatial aggregated point patterns in seismic science <i>Elvira Romano, Jonatan González Monsalve, Francisco Javier Rodríguez Cortés, Jorge Mateu</i>	419
ROC curves with binary multivariate data <i>Lidia Sacchetto, Mauro Gasparini</i>	420
Silhouette-based method for portfolio selection <i>Marco Scaglione, Carmela Iorio, Antonio D'Ambrosio</i>	424
Item weighted Kemeny distance for preference data <i>Mariangela Sciandra, Simona Buscemi, Antonella Plaia</i>	428
A fast and efficient modal EM algorithm for Gaussian mixtures <i>Luca Scrucca</i>	432
Probabilistic archetypal analysis <i>Sohan Seth</i>	436
Multilinear tests of association between networks <i>Daniel K. Sewell</i>	438

Use of multi-state models to maximise information in pressure ulcer prevention trials <i>Linda Sharples, Isabelle Smith, Jane Nixon</i>	442
Partial least squares for compositional canonical correlation <i>Violetta Simonacci Massimo Guarino, Michele Gallo</i>	445
Dynamic modelling of price expectations <i>Rosaria Simone, Domenico Piccolo, Marcella Corduas</i>	449
Towards axioms for hierarchical clustering of measures <i>Philipp Thomann, Ingo Steinwart, Nico Schmid</i>	453
Influence of outliers on cluster correspondence analysis <i>Michel Van de Velden, Alfonso Iodice D'Enza, Lisa Schut</i>	454
Earthquake clustering and centrality measures <i>Elisa Varini, Antonella Peresan, Jiancang Zhuang</i>	458
Co-clustering high dimensional temporal sequences summarized by histograms <i>Rosanna Verde, Antonio Irpino, Antonio Balzanella</i>	462
Statistical analysis of item pre-knowledge in educational tests: latent variable modelling and optimal statistical decision <i>Chen Yunxiao, Lu Yan, Irimi Moustaki</i>	466
Evaluation of the web usability of the University of Cagliari portal: an eye tracking study <i>Gianpaolo Zammarchi, Francesco Mola</i>	468
Application of survival analysis to critical illness insurance data <i>David Zapletal, Lucie Kopecka</i>	472

MODELLING UNOBSERVED HETEROGENEITY OF RANKING DATA WITH THE BAYESIAN MIXTURE OF EXTENDED PLACKETT-LUCE MODELS

Cristina Mollica¹ and Luca Tardella²

¹ Dipartimento di Metodi e Modelli per l'Economia, il Territorio e la Finanza, Sapienza Università di Roma, (e-mail: cristina.mollica@uniroma1.it)

² Dipartimento di Scienze Statistiche, Sapienza Università di Roma, (e-mail: luca.tardella@uniroma1.it)

ABSTRACT: The Plackett-Luce distribution (PL) is one of the most successful parametric options within the class of multistage ranking models to learn the preferences on a given set of items from a sample of ordered sequences. It postulates that the ranking process is carried out by sequentially assigning the positions according to the *forward order*, that is, from the top (most-liked) to the bottom (least-liked) alternative. This assumption has been relaxed with the *Extended Plackett-Luce model* (EPL), thanks to the introduction of the *reference order* parameter describing the rank attribution path. Starting from the recent formulation of the Bayesian EPL, in this work we investigate the further extension into the finite mixture approach as a method to explore the group structure of ranking data.

KEYWORDS: Ranking data, Plackett-Luce model, mixture model, Gibbs sampling, Metropolis-Hastings algorithm.

1 Introduction

A *ranking* is an ordered sequence resulting from the comparative evaluation of a given set of *items* according to a specific criterion. This framework is typical in several areas of research, involving surveys on preferences for consumer goods, psychological/behavioral studies on attitudes, voting systems and the competition/sport context, see Marden, 1995 for a broad review of the statistical literature on methods and models for analysing ranking data.

Formally, a ranking of K items is a vector $\pi = (\pi(1), \dots, \pi(K))$, where the entry $\pi(i)$ indicates the position attributed to the i -th alternative. Data can be equivalently collected in the ordering format $\pi^{-1} = (\pi^{-1}(1), \dots, \pi^{-1}(K))$, where the component $\pi^{-1}(j)$ denotes the item ranked in the j -th position. Thus, ranking data take values in the set of permutations \mathcal{S}_K of the first K integers.

This work concentrates on the parametric family of stagewise models. In particular, our interest is in the *Extended Plackett-Luce model* (EPL), originally proposed by Mollica & Tardella, 2014, both in its basic form and into the finite mixture framework. In that work, inference on the EPL mixture was addressed in the frequentist domain via the EM algorithm. Starting from the recent contribution by Mollica & Tardella, 2019, here we explored the further extension of the Bayesian EPL into the finite mixture approach.

2 The Bayesian EPL mixture

2.1 The Extended Plackett-Luce model

The EPL proposed by Mollica & Tardella, 2014 relies on the relaxation of the conventional forward order assumption of the popular PL class through the introduction of the *reference order* parameter $\rho = (\rho(1), \dots, \rho(K))$, indexing the position assignment order. So, the generic entry $\rho(t)$ indicates the rank attributed at the t -th stage of the ranking process and the entire vector ρ is a discrete parameter represented by a permutation of the first K integers. The probability of a generic ordering under the EPL can be written as

$$\mathbf{P}_{\text{EPL}}(\pi^{-1} | \rho, \underline{p}) = \mathbf{P}_{\text{PL}}(\pi^{-1} \rho | \underline{p}) = \prod_{t=1}^K \frac{p_{\pi^{-1}(\rho(t))}}{\sum_{v=t}^K p_{\pi^{-1}(\rho(v))}} \quad \pi^{-1} \in \mathcal{S}_K.$$

The support parameters p_i 's are proportional to the probabilities for each item to be selected at the first stage and, hence, to be ranked in the position indicated by the first entry of ρ .

2.2 Mixture model setup

In the EPL finite mixture scenario, one assumes that the random sample of N orderings $\underline{\pi}^{-1} = (\pi_1^{-1}, \dots, \pi_N^{-1})$ is drawn from an *heterogenous population* represented by a convex combination of G *subpopulations* (or *groups*), each of which is modelled with a specific EPL distribution. Formally, we set

$$\underline{\pi}_s^{-1} | \underline{\rho}, \underline{p}, \underline{\omega} \stackrel{\text{iid}}{\sim} \sum_{g=1}^G \omega_g \mathbf{P}_{\text{EPL}}(\underline{\pi}_s^{-1} | \rho_g, \underline{p}_g),$$

where ρ_g , \underline{p}_g and ω_g are, respectively, the reference order, the support parameters and the weight of the g -th mixture component.

In order to make Bayesian inference for the G -component EPL mixture analytically tractable, a joint data augmentation strategy combining two sets of latent variables has to be suitably introduced, specifically:

1. the unobserved group labels of each sample unit $s = 1, \dots, N$

$$z_{sg} = \begin{cases} 1 & \text{if unit } s \text{ belongs to the } g\text{-th mixture component,} \\ 0 & \text{otherwise;} \end{cases}$$

2. the latent quantitative variables $\underline{y} = (y_{st})$ for $s = 1, \dots, N$ and $t = 1, \dots, K$, associated to each entry of the data matrix and linked to the component memberships \underline{z} through the following parametric assumption

$$f(\underline{y} | \underline{\pi}^{-1}, \underline{z}, \underline{\rho}, \underline{p}) = \prod_{s=1}^N \prod_{t=1}^K f_{\text{Exp}} \left(y_{st} \mid \prod_{g=1}^G \left(\sum_{v=t}^K p_g \pi_s^{-1}(\rho_g(v)) \right)^{z_{sg}} \right).$$

Thus, the complete-data likelihood can be written as

$$L_c(\underline{\rho}, \underline{p}, \underline{\omega}, \underline{y}, \underline{z}) = \prod_{s=1}^N \prod_{g=1}^G \left(\omega_g \prod_{i=1}^K p_{gi} e^{-p_{gi} \sum_{t=1}^K \delta_{stig} y_{st}} \right)^{z_{sg}},$$

where

$$\delta_{stig} = \begin{cases} 1 & \text{if } i \in \{\pi_s^{-1}(\rho_g(t)), \dots, \pi_s^{-1}(\rho_g(K))\}, \\ 0 & \text{otherwise.} \end{cases}$$

To complete the Bayesian model specification, we considered the following joint prior distribution for the unknown parameters $(\underline{\rho}, \underline{p}, \underline{\omega})$

$$\rho_g \stackrel{\text{iid}}{\sim} \text{Unif}\{\mathcal{S}_K\} \quad p_{gi} \stackrel{i}{\sim} \text{Ga}(c_{gi}, d_g) \quad \underline{\omega} \sim \text{Dir}(\alpha_1, \dots, \alpha_G),$$

where the Gamma densities are indexed by the shape and rate parameters.

2.3 Estimation via MCMC methods

Under the Bayesian model setup described in Section 2.2, the MCMC method proposed by Mollica & Tardella, 2019 to estimate the basic EPL can be easily adapted for the G -component EPL mixture. The outline of the $(l+1)$ -th iteration of the tuned joint Metropolis-Hasting within Gibbs sampling algorithm to

approximate the posterior distribution turns out to be

$$\begin{aligned}
\boldsymbol{\omega}^{(l+1)} | \underline{z}^{(l)} &\sim \text{Dir} \left(\alpha_1 + N_1^{(l)}, \dots, \alpha_G + N_G^{(l)} \right), \\
\rho_g^{(l+1)}, \underline{p}'_g | \underline{\pi}^{-1}, \underline{z}^{(l)} &\sim \mathcal{K}_{\text{TJM}} \circ \mathcal{K}_{\text{SM}}, \\
y_{st}^{(l+1)} | \underline{\pi}_s^{-1}, \underline{z}_s^{(l)}, \underline{\rho}^{(l+1)}, \underline{p}' &\sim \text{Exp} \left(\prod_{g=1}^G \left(\sum_{i=1}^K \delta_{stig}^{(l+1)} p'_{gi} \right)^{z_{sg}^{(l)}} \right), \\
p_{gi}^{(l+1)} | \underline{\pi}^{-1}, y_s^{(l+1)}, \underline{z}^{(l)}, \rho_g^{(l+1)} &\sim \text{Ga} \left(c_{gi} + N_g^{(l)}, d_g + \sum_{s=1}^N z_{sg}^{(l)} \sum_{t=1}^K \delta_{stig}^{(l+1)} y_{st}^{(l+1)} \right), \\
z_s^{(l+1)} | \underline{\pi}_s^{-1}, y_s^{(l+1)}, \underline{\rho}^{(l+1)}, \underline{p}^{(l+1)}, \underline{\omega}^{(l+1)} &\sim \text{Multinom} \left(1, \left(m_{s1}^{(l+1)}, \dots, m_{sG}^{(l+1)} \right) \right),
\end{aligned}$$

where $N_g^{(l)} = \sum_{s=1}^N z_{sg}^{(l)}$ and

$$m_{sg}^{(l+1)} \propto \omega_g^{(l+1)} \prod_{i=1}^K p_{gi}^{(l+1)} e^{-p_{gi}^{(l+1)} \sum_{t=1}^K \delta_{stig}^{(l+1)} y_{st}^{(l+1)}}.$$

With $\mathcal{K}_{\text{TJM}} \circ \mathcal{K}_{\text{SM}}$ we denote the composition of two kernels, namely a tuned joint Metropolis (TJM) and a local swap move (SM) needed to solve the reference order simulation step and ensure an adequate mixing. For the mixture setting, the TJM and the SM are performed on the subsamples determined by the group memberships to iteratively draw the specific reference orders ρ_g .

The determination of the optimal number of mixture components can be addressed with the popular DIC (Spiegelhalter *et al.*, 2002).

References

- MARDEN, J. I. 1995. *Analyzing and modeling rank data*. Monographs on Statistics and Applied Probability, vol. 64. Chapman & Hall.
- MOLLIKA, C., & TARDELLA, L. 2014. Epitope profiling via mixture modeling of ranked data. *Statistics in Medicine*, **33**(21), 3738–3758.
- MOLLIKA, C., & TARDELLA, L. 2017. Bayesian mixture of Plackett-Luce models for partially ranked data. *Psychometrika*, **82**(2), 442–458.
- MOLLIKA, C., & TARDELLA, L. 2019. Bayesian analysis of ranking data with the Extended Plackett-Luce model. (*submitted*).
- SPIEGELHALTER, D. J., BEST, N. G., CARLIN, B. P., & VAN DER LINDE, A. 2002. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **64**(4), 583–639.