

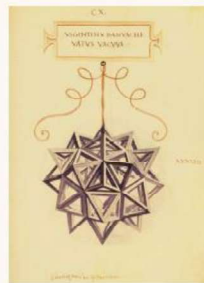
DIPARTIMENTO DI ECONOMIA E GIURISPRUDENZA
UNIVERSITÀ DI CASSINO E DEL LAZIO MERIDIONALE



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

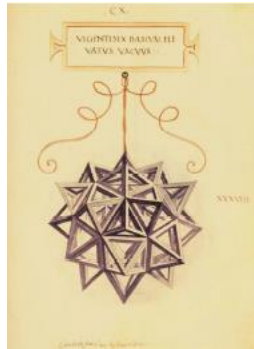


© CC – Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)
<https://creativecommons.org/licenses/by-nc/4.0/>

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

Preface

This book collects the short papers presented at CLADAG 2019, the 12th Scientific Meeting of the Classification and Data Analysis Group (CLADAG) of the Italian Statistical Society (SIS).

The meeting has been organized by the Department of Economics and Law of the University of Cassino and Southern Lazio, under the auspices of the SIS and the International Federation of Classification Societies (IFCS). CLADAG is a member of the IFCS, a federation of national, regional, and linguistically-based classification societies. It is a non-profit, non-political scientific organization, whose aims are to further classification research.

Every two years, CLADAG organizes a scientific meeting, devoted to the presentation of theoretical and applied papers on classification and related methods of data analysis in the broad sense. This includes advanced methodological research in multivariate statistics, mathematical and statistical investigations, survey papers on the state of the art, real case studies, papers on numerical and algorithmic aspects, applications in special fields of interest, and the interface between classification and data science. The conference aims at encouraging the interchange of ideas in the above-mentioned fields of research, as well as the dissemination of new findings.

CLADAG conferences, initiated in 1997 in Pescara (Italy), were soon considered as an attractive information exchange market and became a most important meeting point for people interested in classification and data analysis. One reason was

certainly the fact that a selection of the presented papers is regularly published in (post-conference) proceedings, typically by Springer Verlag.

The Scientific Committee of CLADAG2019 conceived the Plenary and Invited Sessions to provide a fresh perspective on the state of the art of knowledge and research in the field. The scientific program of CLADAG 2019 is particularly rich. All in all, it comprises 5 Keynote Lectures, 32 Invited Sessions promoted by the members of the Scientific Program Committee, 16 Contributed Sessions, a Round Table and a Data Competition. We thank all the session organizers for inviting renowned speakers, coming from 28 countries. We are greatly indebted to the referees, for the time spent in a careful review.

The editors would like to express their gratitude to the Rector of the University of Cassino and Southern Lazio and the Director of the Department of Economics and Law for having hosted the meeting. Special thanks are finally due to the members of the Local Organizing Committee and all the people who with their abnegation and enthusiasm have worked for CLADAG 2019.

Special thanks go to Alfiero Klain and Livia Iannucci for the editorial and administrative support.

Last but not least, we thank all the authors and participants, without whom the conference would not have been possible.

Cassino, September 11, 2019

Giovanni C. Porzio
Francesca Greselin
Simona Balzano

ROBUST MODEL-BASED CLUSTERING WITH MILD AND GROSS OUTLIERS

Alessio Farcomeni¹ and Antonio Punzo²

¹ Department of Public Health and Infectious Diseases, Sapienza - University of Rome,
(e-mail: alessio.farcomeni@uniroma1.it)

² Department of Economics and Business, University of Catania,
(e-mail: antonio.punzo@unict.it)

ABSTRACT: We propose a model-based clustering procedure for mild and gross outliers. Our mixture model is based on heavy-tailed components (e.g., the contaminated normal distribution), but it is assumed to apply only to a subset of the data. Consequently, a proportion of observations is trimmed. We propose a penalized likelihood approach for estimation and selection of the proportions of mild and gross outliers, where the penalty parameter is fixed by formal optimality arguments. We conclude with an original real data example on the identification of the source from illicit drug shipments seized in Italy and Spain.

KEYWORDS: `tclust`, contaminated normal, penalized likelihood.

1 Introduction

In clustering based on the normal mixture model there are two main approaches to deal with contamination. One is based on the use of heavy-tailed or skewed component distributions. A recent example in this direction, preserving elliptical contours of clusters, are mixtures of contaminated normal (CN) distributions (Punzo & McNicholas, 2016). Component-wise methods are well suited to work with mild outliers (Ritter, 2015), and are sometimes labeled as weakly robust. A separate body of literature has instead worked with outliers in more general position, including gross outliers, and has usually proceeded by discarding or at least downweighting a proportion of the observations (Farcomeni & Greco, 2015). A good example is `tclust` (García-Escudero *et al.*, 2008), where a fixed proportion of observations is trimmed and the rest is assumed to follow a normal mixture model. These procedures have often formal robustness properties, e.g., positive breakdown point asymptotically.

In this work we merge the two approaches above by estimating a CN mixture after trimming a fixed proportion of gross outliers. Our model can be

seen from two different perspectives. On the one hand, clusters having a distribution with slightly heavy tails might be desired in order to assign as many observations to clusters as possible. In this case, it is indeed assumed that clean observations arise from, for example, a CN model. On the other hand, the trade off between mild and gross outliers is exploited in order to increase efficiency: some (mild) outliers are assigned to a cluster and contribute to centroid estimation, therefore decreasing the final mean squared error (MSE).

In this work we tackle also an additional open problem with trimming procedures, that of selecting the trimming proportion. Our proposal is based on a penalized likelihood approach, where the trimming proportion is in practice substituted by a penalty parameter. The advantage is that we can identify a heuristic but theoretically justified way of choosing an optimal penalty level, and therefore an optimal trimming proportion. Our fixed-penalty approach in some sense solves the issue of selecting the trimming proportion both for our model and the special case of trimmed normal mixture models (`tclust`). The methodology proposed in this paper has been implemented in R functions which can be downloaded from <https://github.com/afarcome/cntclust>.

2 Methodology

Let $x_1, \dots, x_i, \dots, x_n$ be a sample of n observations in d dimensions. Moreover, let $\alpha_0 \geq 0$ denote a trimming proportion of outliers which shall not be used to estimate model parameters. We assume data arise from the contaminated spurious outlier model

$$\prod_{i \in R} \sum_{j=1}^k \pi_j f_{\text{CN}}(x_i; \mu_j, \Sigma_j, \alpha_j, \eta_j) \prod_{i \notin R} g_i(x_i), \quad (1)$$

where R denotes a set of non-trimmed observations of cardinality $\lfloor (1 - \alpha_0)n \rfloor$ and g_i are pdfs generating the outliers in general position. Let $f_{\text{N}}(\cdot; \mu, \Sigma)$ denote the probability density function (pdf) of a d -variate normal (N) distribution with mean vector μ and covariance matrix Σ . In (1), $f_{\text{CN}}(x; \mu, \Sigma, \alpha, \eta) = (1 - \alpha) f_{\text{N}}(x; \mu, \Sigma) + \alpha f_{\text{N}}(x; \mu, \eta \Sigma)$ denotes the pdf of a d -variate CN distribution with mean vector μ , scale matrix Σ , proportion of mild outliers $\alpha \in (0, 1)$, and degree of contamination $\eta > 1$.

To estimate the parameters, we optimize the profile likelihood

$$\ell(\vartheta) = \sum_{j=1}^k \sum_{i \in R_j} \ell_i(\vartheta) = \sum_{j=1}^k \sum_{i \in R_j} [\ln \pi_j + \ln f_{\text{CN}}(x_i; \mu_j, \Sigma_j, \alpha_j, \eta_j)], \quad (2)$$

where R_j denotes the set of observations assigned to the j -th cluster. To make maximization of (2) a well defined problem, we adopt the classical eigenvalue ratio constraint proposed by García-Escudero *et al.*, 2008.

Model (1) involves the difficult choice of $\alpha_0, \alpha_1, \dots, \alpha_k$, where α_0 controls the proportion of gross outliers and α_j the proportion of mild outliers in the j -th cluster. We propose a LASSO-type penalized likelihood approach enforcing a sparse model selection in which some values in the set $(\alpha_0, \alpha_1, \dots, \alpha_k)$ might be set to zero. A general form of penalized log-likelihood is given by

$$\ell(\vartheta) + P(\alpha_0, \alpha_1, \dots, \alpha_k), \quad (3)$$

and we propose using $P(\alpha_0, \dots, \alpha_k) = -\log(n) \sum_{j=0}^k v_j \alpha_j$. In order to reduce the number of penalty parameters, we set $v_0 = nv$ and $v_j = v$ for $j > 0$.

The choice of the penalty parameter v has got direct consequences on the estimated trimming proportion α_0 . If also $\alpha_1, \dots, \alpha_k$ are included in the penalty, it also affects their estimates. Surprisingly enough, mapping the problem of selecting contaminating proportions to the scale of the likelihood gives an asymptotically “optimal” fixed value, $v = \sqrt{2d}$, which under certain assumptions guarantees that observations outside a chi-square type ellipse from a bulk of the data are trimmed.

Maximization of (2), and for fixed v of (3), is carried out using a classification expectation-conditional maximization (CECM) algorithm, where eigenvalue ratio constraints are activated at the conditional maximization step is needed,

3 Example about clustering illicit drug shipments

We analyze data about $n = 151$ seizures of shipments of cocaine and heroin in Italy and Spain. They were sent to the forensic laboratories for checking the nature of the substance and quantifying the absolute and relative contents of each of several chemical compounds. In modern forensics it is believed that the contents of certain solvents might be useful for identifying the source, that is, clustering packages with respect to the illicit laboratory where the drug was processed. We verify this assumption by focusing on $d = 3$ compounds: hexane, acetone, and 2-propanol. We fix $k = 2$ and estimate a classical normal mixture model and a contaminated normal mixture model without trimming first. Then we use robust clustering methods: `tclust` and the contaminated-normal mixture model with trimming. In Table 1 we report, for values of the trimming level chosen using our penalized likelihood approach, the adjusted

Rand-index (ARI) showing the agreement between the class labels and the true underlying Italy/Spain location of seizure. With no or insufficient trimming one might conclude that there is no relationship between solvent contents and seizure location. On the other hand, after trimming the agreement becomes fairly high. As expected we note that the optimal trimming level using `tclust` is slightly larger than those using `CNTCLUST0`. While in our low sample size example this might not have strong consequences in terms of MSE, $\lceil 151(0.066 - 0.053) \rceil = 2$ seizures will not be attributed to a location using `tclust`, which can have forensic consequences.

Table 1. Adjusted Rand-index (ARI) for location of drug seizure and clustering. In parentheses the trimming level. NM: normal mixture, CNM: contaminated normal mixture, `tclust`: trimmed NM, `CNTCLUST0`: trimmed CNM, `CNTCLUST`: penalized trimmed CNM with $v = \sqrt{2d}$ and fixed trimming level. The trimming level selected with our fixed-penalty approach is indicated with $\hat{\alpha}_0$.

Method	ARI	Method	ARI
NM	-0.076	<code>CNTCLUST0</code> ($\hat{\alpha}_0 = 0.053$)	0.660
CNM	-0.069	<code>CNTCLUST</code> ($\hat{\alpha}_0 = 0.053$)	0.658
<code>tclust</code> ($\hat{\alpha}_0 = 0.066$)	0.657		

References

- FARCOMENI, A., & GRECO, L. 2015. *Robust Methods for Data Reduction*. Boca Raton, FL: CRC Press.
- GARCÍA-ESCUADERO, L. A., GORDALIZA, A., MATRAN, C., & MAYO-ISCAR, A. 2008. A general trimming approach to robust cluster analysis. *Annals of Statistics*, **36**, 1324–1345.
- PUNZO, A., & MCNICHOLAS, P. D. 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, **58**(6), 1506–1537.
- RITTER, G. 2015. *Robust Cluster Analysis and Variable Selection*. CRC Press.