







CLADAG 2013

9thSCIENTIFIC MEETING OF THE CLASSIFICATION AND DATA ANALYSIS GROUP OF THE ITALIAN STATISTICAL SOCIETY

September 18 - 20, 2013

University of Modena and Reggio Emilia San Geminiano Complex - Modena, Italy

Book of Abstracts

Editors: Tommaso Minerva, Isabella Morlini, Francesco Palumbo CLEUP ISBN: 9788867871179

Patronage





Table of Contents

Tommaso Agasisti, Patrizia Falzetti Socioeconomic sorting and test scores:an empirical analysis in the Italian junior secondary schools	pag. 2
Dario Albarello, Vera D'Amico Empirical testing of probabilistic seismic hazardmodels	pag. 9
Federico Andreis, Pier Alda Ferrari A proposal for the multidimensional extension of CUB models	pag. 15
Morten Arendt Rasmussen, Evrim Acar Data fusion in the framework of coupled matrix tensor factorization withcommon, partially common and unique factors	pag. 19
Luigi Augugliaro, Angelo M. Mineo Estimation of Sparse Generalized LinearModels: the dglars package	pag. 20
Antonio Balzanella, Lidia Rivoli, Elvira Romano A comparison between two tools for data streamsummarization	pag. 24
Lucio Barabesi, Giancarlo Diana, Pier Francesco Perri Gini Index Estimation in Randomized ResponseSurveys	pag. 28
Francesco Bartolucci, Federico Belotti, Franco Peracchi A test for time-invariant individual effects ingeneralized linear models for panel data	pag. 32
Erich Battistin, Carlos Lamarche, Enrico Rettore <i>Identification of the distribution of the causal effect of an intervention</i> <i>using a generalised factor model</i>	pag. 36
Matilde Bini, Lucio Masserini Internal effectiveness of educational offer andstudents' satisfaction: a SEM approach	pag. 37
Matilde Bini, Leopoldo Nascia, Alessandro Zeli Groups heterogeneity and sectorsconcentration: a structural equation modelingfor micro level analysis of firms	pag. 41
Giuseppe Boari, Marta Nai Ruscone Use of Relevant Principal Components to Definea Simplified Multivarate Test Procedure ofOptimal Clutering	pag. 45
Giuseppe Boari, Gabriele Cantaluppi, Angelo Zanella Some Distance Proposals for Cluster Analysis inPresence of Ordinal Variables	pag. 49

Laura Bocci, Donatella Vicari A general model for INDCLUS with externalinformation	pag. 53
Paola Bongini, Paolo Trivellato, Mariangela Zenga <i>The financial literacy and the undergraduates</i>	pag. 57
Riccardo Bramante, Marta Nai Ruscone, Pasquale Spani Credit risk measurement and ethical issue: someevidences from the italian banks	pag. 61
Pierpaolo Brutti, Lucio Ceccarelli, Fulvio De Santis, Stefania Gubbiotti On the Stylometric Authorship of Ovid's DoubleHeroides: An Ensemble Clustering Approach	pag. 65
Silvia Caligaris, Fulvia Mecatti and Patrizia Farina Causal Inference in Gender Discrimination inChina: Nutrition, Health, Care	pag. 69
Giorgio Calzolari, Antonino Di Pino Self-Selection and Direct Estimation of Across-Regime Correlation Parameter	pag. 73
Maria Gabriella Campolo, Antonino Di Pino, Ester Lucia Rizzi Modern Vs. Traditional: A cluster-basedspecification of gender and familistic attitudesand their influence on the division of labour of Italia couples	pag. 77 n
Gabriele Cantaluppi, Marco Passarotti Clustering the Four Gospels in the Greek, Latin, Gothic and Old Church SlavonicTranslations	pag. 81
Carmela Cappelli, Francesca Di Iorio Regression Trees for change point analysis:methods, applications and recent developments	pag. 85
Roberto Casarin and Marco Tronzano and Domenico Sartore <i>Bayesian Stochastic Correlation Models</i>	pag. 89
Rosalia Castellano, Gennaro Punzo, Antonella Rocca <i>Evaluating the selection effect in labour marketswith a low female</i> <i>participation</i>	pag. 93
Paola Cerchiello, Paolo Giudici A statistical based H index for the evaluation ofe-markets	pag. 97
Annalisa Cerquetti Bayesian nonparametric estimation of globaldisclosure risk	pag. 101
Enrico Ciavolino, Roberto Savona <i>The Forecasting side of Sovereign Risk: aGeneralized Cross Entropy Ap</i>	pag. 105 oproach

Nicoletta Cibella, Tiziana Tuoto, Luca Valentino What data tell you that models can't say	pag. 109
Roberto Colombi, Sabrina Giordano Multiple Hidden Markov Models for CategoricalTime Series	pag. 114
Pier Luigi Conti, Daniela Marella Asymptotics in survey sampling for high entropysampling designs	pag. 118
Claudio Conversano, Massimo Cannas, Francessco Mola On the Use of Recursive Partitioning in CasualInference: A Proposal	pag. 122
Franca Crippa, Marcella Mazzoleni, Mariangela Zenga <i>Keeping the pace with higher education. A fuzzystates gender study</i>	pag. 128
F. Cugnata, C. Guglielmetti and S. Salini CUB model to validate FACIT TS-PSmeasurement instrument	pag. 133
Rosario D'Agata, Venera Tomaselli Multilevel Approach in Meta-Analysis of Pre-Election Poll Accuracy	pag. 137
Alfonso Iodice D'Enza and Angelos Markos Low-dimensional tracking of associationstructures in categorical data	pag. 141
Giulio D'Epifani Self-censored Categorical ResponsesA device for recovering latent behaviors	pag. 145
Pierpaolo D'Urso, Marta Disegna, Riccardo Massari Tourism Market Segmentation with ImpreciseInformation	pag. 150
Utkarsh J. Dang, Salvatore Ingrassia, Paul D. McNicholas and Ryan Browne <i>Cluster-weighted models for multivariateresponse and extensions</i>	pag. 154
Cristina Davino, Domenico Vistocco Unsupervised Classification through QuantileRegression	pag. 158
F. Marta L. Di Lascio, Simone Giannerini A copula-based approach to discoverinter-cluster dependence relationships	pag. 162
Josè G. Dias, Sofia B. Ramos Hierarchical market structure of Euro arearegime dynamics	pag. 166
Drago Carlo, Balzanella Antonio Consensus Community Detection: a NonmetricMDS Approach	pag. 170
Fabrizio Durante, Roberta Pappad`a and Nicola Torelli Clustering financial time series by measures oftail dependence	pag. 174

Marco Enea, Antonella Plaia Influence diagnostics for generalized linearmixed models: a gradient-like statistic	pag. 178
Enrico Fabrizi, Maria R. Ferrante, Carlo Trivisano Joint estimation of poverty and inequalityparameters in small areas	pag. 182
Giorgio Fagiolo, Andrea Roventini Macroeconomic Policy in DSGE andAgent-Based Models	pag. 187
Salvatore Fasola, Mariangela Sciandra New Flexible Probability Distributions forRanking Data	pag. 191
Maria Brigida Ferraro, Paolo Giordani A new fuzzy clustering algorithm with entropyregularization	pag. 195
Camilla Ferretti, Piero Ganugi, Renato Pieri Mobility measures for the dairy farms inLombardy	pag. 199
Silvia Figini, Marika Vezzoli Model averaging and ensemble methods for riskcorporate estimation	pag. 203
Luis Angel García-Escudero, Alfonso Gordaliza, Carlos Matrán, Agustín Mayo-Iscar New proposals for clustering based on trimmingand restrictions	pag. 207
Andreas Geyer-Schulz, Fabian Ball Formal Diagnostics for Graph Clustering: TheRole of Graph Automorphisms	pag. 211
Massimiliano Giacalone, Angela Alibrandi An overview on multiple regression models basedon permutation tests	pag. 215
Francesca Giambona, Mariano Porcu <i>The determinants of Italian students' readingscores:</i> <i>a Quantile Regression analysis</i>	pag. 219
Paolo Giordani, Henk A.L. Kiers, Maria Antonietta Del Ferraro <i>The R Package ThreeWay</i>	pag. 223
Giuseppe Giordano, Ilaria Primerano Co-occurence Network from SemanticDifferential Data	pag. 227
Paolo Giudici Financial risk data analysis	pag. 231
Silvia Golia, Anna Simonetto A Comparison between SEM and Rasch model:the polytomous case	pag. 237
Anna Gottard Some considerations on VCUB models	pag. 241

Francesca Greselin, Salvatore Ingrassia Data driven EM constraints for mixtures offactor analyzers	pag. 245
Leonardo Grilli, Carla Rampichini, Roberta Varriale <i>Predicting students' academic performance: achallenging issue in</i> <i>statistical modelling</i>	pag. 249
Luigi Grossi, Fany Nan Robust estimation of regime switching models	pag. 255
Kristian Hovde Liland Variable selection in sequential multi-block analysis	pag. 259
Maria Iannario Robustness issues for a class of models forordinal data	pag. 260
Maria Iannario, Domenico Piccolo A class of ordinal data models in R	pag. 264
Salvatore Ingrassia, Antonio Punzo Parsimony in Mixtures with Random Covariates	pag. 268
Hiroshi Inoue International Relations Based on the VotingBehavior in General Assembly	pag. 272
Carmela Iorio, Massimo Aria, Antonio D'Ambrosio Visual model representation and selection forclassification and regression trees	pag. 276
Monia Lupparelli, Luca La Rocca, Alberto Roverato Log-Mean Linear Parameterizations for SmoothIndependence Models	pag. 284
Marica Manisera, Paola Zuccolotto Nonlinear CUB models	pag. 288
Marica Manisera, Marika Vezzoli Finding number of groups using a penalizedinternal cluster quality index	pag. 292
Daniela Marella, Paola Vicard <i>Object-Oriented Bayesian Network to deal withmeasurement error</i> <i>in household surveys</i>	pag. 296
Angelos Markos, Alfonso Iodice D'Enza, Michel Van de Velden Beyond tandem analysis: joint dimensionreduction and clustering in R	pag. 300
F. Martella and M. Alfò A biclustering approach for discrete outcomes	pag. 304

Mariagiulia Matteucci, Stefania Mignani, Roberto Ricci A Multidimensional IRT approach to analyzelearning achievement of Italian students	pag. 309
Sabina Mazza Extending the Forward Search to theCombination of Multiple Classifiers: A Proposal	pag. 314
Fulvia Mecatti, M. Giovanna Ranalli <i>Plug-in Bootstrap for Sample Survey Data</i>	pag. 318
Alessandra Menafoglio, Matilde Dalla Rosa and Piercesare Secchi A BLU Predictor for Spatially DependentFunctional Data of a Hilbert Space	pag. 322
Maria Adele Milioli, Lara Berzieri, Sergio Zani Comparing fuzzy and multidimensional methodsto evaluate well-being at regional level	pag. 326
Michelangelo Misuraca, Maria Spano Comparing text clustering algorithms from amultivariate perspective	pag. 331
Cristina Mollica, Luca Tardella Mixture models for ranked data classification	pag. 335
Isabella Morlini, Stefano Orlandini Cluster analysis of three-way atmospheric data	pag. 339
Roberto Nardecchia, Roberto Sanzo, Margherita Velucchi, Alessandro Zeli Productivity transition probabilities: A microlevel data analysis for Italian manufacturingsectors (1998-2007)	pag. 345
Andrea Neri, Giuseppe Ilardi Interviewers, co-operation and data accuracy: isthere a link?	pag. 349
Akinori Okada, Satoru Yokoyama Nonhierarchical Asymmetric Cluster Analysis	pag. 353
Marco Perone Pacifico SuRF: Subspace Ridge Finder	pag. 357
Andrea Pagano, Francesca Torti, Jessica Cariboni, Domenico Perrotta <i>Robust clustering of EU banking data</i>	pag. 361
Giuseppe Pandolfo, Giovanni C. Porzio On depth functions for directional data	pag. 365
Andrea Pastore, Stefano F. Tonellato A generalised Silhouette-width measure	pag. 369

Fulvia Pennoni, Giorgio Vittadini Hospital efficiency under two competing paneldata models	pag. 373
Alessia Pini, Simone Vantini The Interval-Wise Control of the Family-WiseError Rate for Testing Functional Data	pag. 377
Mariano Porcu, Isabella Sulis Detecting differences between primary schools inmathematics and reading achievement by usingschools added-value measures of performance	pag. 381
Antonio Punzo, Paul D. McNicholas, Katherine Morris, Ryan P. Browne Outlier Detection via Contaminated MixtureDistributions	pag. 387
Emanuela Raffinetti, Pier Alda Ferrari New perspectives for the RDI index in socialresearch fields	pag. 392
Monia Ranalli, Roberto Rocci <i>Mixture models for ordinal data: a pairwiselikelihood approach</i>	pag. 396
Marco Riani, Andrea Cerioli, Gianluca Morelli Issues in robust clustering	pag. 400
Stèphane Robin Deciphering and modeling heterogeneity ininteraction networks	pag. 404
Rosaria Romano, Francesco Palumbo Partial Possibilistic Regression Path Modeling	pag. 409
Renata Rotondi Classsification of composite seismogenic sourcesthrough probabilitic score indices	pag. 413
Gabriella Schoier On Wild Bootstrap and M Unit Root Test	pag. 417
Luca Scrucca On the implementation of a parallel algorithmfor variable selection in model-based clustering	pag. 421
Paolo Sestito <i>The Role of Learning Measurement in theGovernance of an</i> <i>Education System: anOverview of the Issues</i>	pag. 425
John Shawe-Taylor, Blaz Zlicar Novelty Detection with Support Vector Machines	pag. 430
Nadia Solaro <i>Multidimensional scaling with incompletedistance matrices:</i> <i>an insight into the problem</i>	pag. 431

Luigi Spezia, Cecilia Pinto Markov switching models for high-frequencytime series: flapper skate's depth profile as a casestudy	pag.	435
Ralf Stecking, Klaus B. Schebesch Data Privacy in Credit Scoring: Evaluating SVMApproaches Based on Microaggregated Data	pag.	439
Isabella Sulis, Francesca Giambona, Nicola Tedesco Analyzing university students' careers usingMulti-State Models	pag.	443
Luca Tardella, Danilo Alunni Fegatelli BBRecap for Bayesian BehaviouralCapture-Recapture Modeling	pag.	447
Cristina Tortora, Paul D. McNicholas, Ryan P. Browne <i>Mixtures of generalized hyperbolic factoranalyzers</i>	pag.	451
Giovanni Trovato <i>Testing for endogeneity and countryheterogeneity</i>	pag.	455
Joaquin Vanschoren and Mikio L. Braun, Cheng Soon Ong Open science in machine learning	pag.	461
Valerio Veglio Logistic Regression and Decision Tree:Performance Comparisons in EstimatingCustomers' Risk of Churn	pag.	465
Maurizio Vichi Robust Two-mode clustering	pag.	469
Vincenzina Vitale Hierarchical Graphical Models and ItemResponse Theory	pag.	470
Sara Viviani Extending the JM libraRy	pag.	474
Adalbert F.X. Wilhelm Visualisations of Classification Tree Models: AnEvaluative Comparison	pag.	478

On the Stylometric Authorship of Ovid's Double Heroides: An Ensemble Clustering Approach

Pierpaolo Brutti, Lucio Ceccarelli, Fulvio De Santis, Stefania Gubbiotti

Abstract *Double Heroides* are six elegies traditionally attributed to Ovid, whose authenticity have been repeatedly questioned. As a contribution to establish the period of composition of these elegies, this article proposes a statistical analysis based on consensus clustering of mixed data composed of standard (i.e. scalar) and compositional variables derived by an extensive metrical study of the Ovid's poetic production.

Key words: Consensus clustering, compositional data, Latin metric.

1 Introduction: Ovid's double letters

A relevant open problem in Latin literature is the authorship and the date of Ovids *Heroides XVI-XXI (double Heroides)*. The *Heroides* are a collection of 21 elegies written in *elegiac couplets*. The authenticity of the collection has been questioned since Lackmann [1] who, on the basis of certain metrical anomalies with respect to genuine Ovidian poems as well as of non metrical considerations, claimed that some of the *Heroides* (including the double Heroides) where not composed by Ovid. In this article we propose a statistical approach to provide further elements of discussion on the dating problem. In particular our attempt is to find statistical support, based on metrical features, to the hypotheses that *double Heroides* belong to one of the three main phases of Ovid's poetic production: the early Roman period, the mature Roman period and the exile period. The stylometric methodology we adopt is based on consensus clustering techniques for mixed data composed of standard (i.e. scalar) and compositional variables.

Lucio Ceccarelli

Università degli Studi dell'Aquila e-mail: lucio.ceccarelli@univaq.it

Pierpaolo Brutti, Fulvio De Santis, Stefania Gubbiotti Sapienza Università di Roma, e-mail: stefania.gubbiotti@uniroma1.it

2 Basics of metric: the elegiac couplet

Meter is the basic rhythmic structure of a verse. The study of meters and forms of versification is known as metric. We here focus on the so called *elegiac couplet* or *distich*, a poetic form initially introduced in Greek lyric, later on adopted by Roman poets, and in particular by Ovid. To illustrate the "anatomy" of the elegiac couplet, let us consider one of the 406 distichs from Ovid's *Fasti VI*: The main steps in the

Tempora labuntur, tacitisque senescimus annis;

et fugiunt freno non remorante dies.

metrical analysis are summarized as follows (see Figure 2).

- Each syllable of the words of a verse is categorized as *long* (−) or *short* (∪) according to its *weight*. (e.g. in Tempora, Tem- is long, -po- and -ra are short).
- Specific sequences of syllables define a *foot* (delimited by |...|), for instance

Dactyl (D)	formed by 3 syllables	-UU
Spondee (S)	formed by 2 syllables	
Trochee (T)	formed by 2 syllables	$-\cup$

- Each verse is formed by a specific number of *feet*. The first verse of the elegiac distich is called *hexameter*; the second one is called *pentameter*.
- The metrical pattern of feet in a verse is then summarized by a sequence of letters. For instance, in the couplet above the scheme is DSDDDS · DD-DD-



Fig. 1 "Anatomy" of the elegiac couplet; (-) and (\cup) denote long and short syllables respectively.

In general, in the hexameter each of the first four feet can be alternatively a *dactyl* or a *spondee*. The fifth foot is almost always a *dactyl*. The sixth foot is either a *spondee* or a *trochee*. Conversely, the pentamer is made up of two equal parts containing two dactyls followed by a long syllable (*hemistich*). Spondees replace dactyls in the first half, but never in the second. The choice of a particular scheme for the verses, together with many other metrical features that will be mentioned in the following section, yield a great variability in the realization of the elegiac distich, that strongly characterizes the style of each single author. For instance, if we consider the first four feet of the hexameter only, there are $2^4 = 16$ possible alternative choices of dactyls and spondees. In summary, the metrical technique is

a very personal ability of the poet and it reflects not only his skills, but also his sensitivity. In this sense, the stylometric analysis can be helpful in the attempt of attributing a poem to a specific author.

3 Stylometric features

The goal of a metric study is to identify characterizing stylistic features of a poet with respect to the metric language of the tradition, as well as to detect deviations of metric features of some parts of his own poetic production with respect to his entire work. This kind of study requires a translation of metrical phenomena into quantitative data, whose relevance can be pointed out uniquely by appropriate statistical analysis of poetic *corpora* (see [2]). All the poetic production in elegiac couplets attributed to Ovid has been examined from a metrical point of view in [3]. In this metrical study many quantitative informations has been extracted and collected on each section of every poem, such as the total frequency of the dactyls and the distribution of dactyls both in the first four feet of the hexameter and in the first half of the pentameter, the presence of short syllables in the couplet, the occurrence of some particular forms (e.g. synalepha, clausula), and so on. Consequently, our dataset consists of 15 distinct metrical features measured on each of the 27 poems considered as statistical units. Most of these 15 variables are actually frequency distributions summarizing the occurrence of a specific metrical phenomenon over a whole poem, and therefore they will be treated as *compositional data* (see [1] and [8]). For example, in Figure 3 are represented three variables observed on two of



Fig. 2 A snippet of the stylometric dataset.

the poems: the first and the second one consist of distributions associated to the realizations of the pentameter scheme, and to the clausula in the exameter respectively, whereas the third scalar one is a particular stylometric quantity usually called *replication index*.

4 Methods and results

Consensus clustering, emerged as an important elaboration of the classical clustering problem. These methods are commonly used to establish consensus among multiple clustering algorithms, or multiple realizations of the same clustering algorithm on a single dataset, or even to integrate multi–source data (for a survey see [7]). In this work we adapt to the peculiarities of the data at hand, two recently proposed tecniques, namely the *Multiple Dataset Integration* method developed in [4] and the Bayesian consensus clustering based on finite Dirichlet mixture models described in [6]. Both approaches use a statistical framework to cluster each data source separately while simultaneously modeling dependence between the clusterings in order to borrow strength across data sources. Aggregation mechanisms of this type make the resulting overall clustering more robust and stable of other unsupervised classification solutions, while still allowing to better interpret the contribution of each data–source to the final partition.

In summary, the consensus clustering shows a quite remarkable distinction between three groups, characterized by different metric profiles. It can be noticed that in each group one of the three periods prevails. Finally, the *double Heroides* turn out to be compatible with the (stylo)metric features of the latest period of Ovid's poetry.

References

- 1. Aitchison J.: The Statistical Analysis of Compositional Data. The Blackburn Press, (2003)
- 2. Ceccarelli L.: Contributi sulla storia dell'esametro latino. Herder, Roma (2008)
- Ceccarelli L.: L'evoluzione del distico elegiaco tra Catullo e Ovidio. In: Cristofoli R., Santini C., Santucci F. Properzio tra tradizione e innovazione, Atti del convegno internazionale (Assisi-Spello, 21-23 maggio 2010), Assisi 2012, pp. 47–97.
- Kirk P., Griffin J.E., Savage R.S., Ghahramani, Z., Wild, D.L.: Bayesian correlated clustering to integrate multiple datasets. Bioinformatics, 28(24), 3290-3297, (2012)
- 5. Lachmann K.: De Ovidi epistulis. Progr. Univers. Berolinensis, (1848)
- 6. Lock E.F., Dunson D.B.: Bayesian consensus clustering. arXiv:1302.7280, (2013)
- Nguyen, N., Caruana, R. Consensus clusterings. In: Proceedings of the 7th IEEE International Conference on Data Mining (ICDM 2007), October 28-31, 2007, Omaha, Nebraska, USA, pp. 607-612. IEEE Computer Society, 2007
- van den Boogaart K.G., Tolosana-Delgado R.: Analyzing Compositional Data with R. Springer, (2013)

4