

# Book of Short Papers SIS 2018

**Editors: Antonino Abbruzzo - Eugenio Brentari**

**Marcello Chiodi - Davide Piacentino**

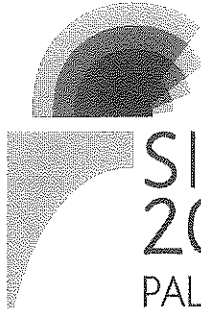


Copyright © 2018

PUBLISHED BY PEARSON

WWW.PEARSON.COM

*First printing, November 2018, ISBN-9788891910233*



**SIS**  
**2018**  
**PALERMO**

49TH SCIENTIFIC MEETING  
OF THE ITALIAN  
STATISTICAL SOCIETY

20-22 JUNE

## Contents

<b>1</b>	<b>Preface</b> .....	<b>v</b>
<b>2</b>	<b>Plenary Sessions</b> .....	<b>19</b>
<b>2.1</b>	<b>A new paradigm for rating data models. <i>Domenico Piccolo</i></b>	<b>19</b>
<b>2.2</b>	<b>Statistical challenges and opportunities in modelling coupled behaviour-disease dynamics of vaccine refusal. <i>Plenary/Chris T. Bauch</i></b>	<b>32</b>
<b>3</b>	<b>Specialized Sessions</b> .....	<b>45</b>
<b>3.1</b>	<b>3.1 - Bayesian Nonparametric Learning</b>	<b>45</b>
3.1.1	Bayesian nonparametric covariate driven clustering. <i>Raffaele Argiento, Ilaria Bianchini, Alessandra Guglielmi and Ettore Lanzarone</i> .....	46
3.1.2	A Comparative overview of Bayesian nonparametric estimation of the size of a population. <i>Luca Tardella and Danilo Alunni Fegatelli</i> .....	56
3.1.3	Logit stick-breaking priors for partially exchangeable count data. <i>Tommaso Rigon</i>	64
<b>3.2</b>	<b>BDsports - Statistics in Sports</b>	<b>72</b>
3.2.1	A paired comparison model for the analysis of on-field variables in football matches. <i>Gunther Schaubberger and Andreas Groll</i> .....	72
3.2.2	Are the shots predictive for the football results?. <i>Leonardo Egidi, Francesco Pauli, Nicola Torelli</i> .....	81
3.2.3	Zero-inflated ordinal data models with application to sport (in)activity. <i>Maria Iannario and Rosaria Simone</i> .....	89
<b>3.3</b>	<b>Being young and becoming adult in the Third Millennium: definition issues and processes analysis</b>	<b>97</b>
3.3.1	Do Social Media Data predict changes in young adults' employment status? Evidence from Italy. <i>Andrea Bonanomi and Emiliano Sironi</i> .....	97

3.3.2	Parenthood: an advanced step in the transition to adulthood. <i>Cinzia Castagnaro, Antonella Guarneri and Eleonora Meli</i> . . . . .	106
<b>3.4</b>	<b>Economic Statistics and Big Data</b>	<b>114</b>
3.4.1	Improvements in Italian CPI/HICP deriving from the use of scanner data. <i>Alessandro Brunetti, Stefania Fatello, Federico Polidoro, Antonella Simone</i> . . . . .	114
3.4.2	Big data and spatial price comparisons of consumer prices. <i>Tiziana Laureti and Federico Polidoro</i> . . . . .	123
<b>3.5</b>	<b>Financial Time Series Analysis</b>	<b>131</b>
3.5.1	Dynamic component models for forecasting trading volumes. <i>Antonio Naimoli and Giuseppe Storti</i> . . . . .	131
3.5.2	Conditional Quantile-Located VaR. <i>Giovanni Bonaccollo, Massimiliano Caporin and Sandra Paterlini</i> . . . . .	140
<b>3.6</b>	<b>Forensic Statistics</b>	<b>146</b>
3.6.1	Cause of effects: an important evaluation in Forensic Science. <i>Fabio Corradi and Monica Musio</i> . . . . .	146
3.6.2	Evaluation and reporting of scientific evidence: the impact of partial probability assignments. <i>Silvia Bozza, Alex Biedermann, Franco Taroni</i> . . . . .	155
<b>3.7</b>	<b>Missing Data Handling in Complex Models</b>	<b>161</b>
3.7.1	Dependence and sensitivity in regression models for longitudinal responses subject to dropout. <i>Marco Alfo' and Maria Francesca Marino</i> . . . . .	161
3.7.2	Multilevel analysis of student ratings with missing level-two covariates: a comparison of imputation techniques. <i>Maria Francesca Marino e Carla Rampichini</i> . . . . .	170
3.7.3	Multilevel Multiple Imputation in presence of interactions, non-linearities and random slopes. <i>Matteo Quartagno and James R. Carpenter</i> . . . . .	175
<b>3.8</b>	<b>Monitoring Education Systems. Insights from Large Scale Assessment Surveys</b>	<b>183</b>
3.8.1	Educational Achievement of Immigrant Students. A Cross-National Comparison Over-Time Using OECD-PISA Data. <i>Mariano Porcu</i> . . . . .	183
<b>3.9</b>	<b>New Perspectives in Time Series Analysis</b>	<b>192</b>
3.9.1	Generalized periodic autoregressive models for trend and seasonality varying time series. <i>Francesco Battaglia and Domenico Cucina and Manuel Rizzo</i> . . . . .	192
<b>3.10</b>	<b>Recent Advances in Model-based Clustering</b>	<b>201</b>
3.10.1	Flexible clustering methods for high-dimensional data sets. <i>Cristina Tortora and Paul D. McNicholas</i> . . . . .	201
3.10.2	A Comparison of Model-Based and Fuzzy Clustering Methods. <i>Marco Alfo', Maria Brigida Ferraro, Paolo Giordani, Luca Scrucca, and Alessio Serafini</i> . . . . .	208
3.10.3	Covariate measurement error in generalized linear models for longitudinal data: a latent Markov approach. <i>Roberto Di Mari, Antonio Punzo, and Antonello Maruotti</i>	216
<b>3.11</b>	<b>Statistical Modelling</b>	<b>224</b>
3.11.1	A regularized estimation approach for the three-parameter logistic model. <i>Michela Battauz and Ruggero Bellio</i> . . . . .	224
3.11.2	Statistical modelling and GAMLSS. <i>Mikis D. Stasinopoulos and Robert A. Rigby and Fernanda De Bastiani</i> . . . . .	233
<b>3.12</b>	<b>Young Contributions to Statistical Learning</b>	<b>239</b>
3.12.1	Introducing spatio-temporal dependence in clustering: from a parametric to a nonparametric approach . <i>Clara Grazian, Gianluca Mastrantonio and Enrico Bibbona</i> . . . . .	239

3.12.2	Bayesian inference for hidden Markov models via duality and approximate filtering distributions. <i>Guillaume Kon Kam King, Omiros Papaspiiliopoulos and Matteo Ruggiero</i> . . . . .	248
3.12.3	K-means seeding via MUS algorithm. <i>Leonardo Egidi, Roberta Pappada, Francesco Pauli, Nicola Torelli</i> . . . . .	256

#### **4 Sollicited Sessions** . . . . . 263

##### **4.1 Advances in Discrete Latent Variable Modelling** **263**

4.1.1	A joint model for longitudinal and survival data based on a continuous-time latent Markov model. <i>Alessio Farcomeni and Francesco Bartolucci</i> . . . . .	264
4.1.2	Modelling the latent class structure of multiple Likert items: a paired comparison approach. <i>Brian Francis</i> . . . . .	273
4.1.3	Dealing with reciprocity in dynamic stochastic block models. <i>Francesco Bartolucci, Maria Francesca Marino, Silvia Pandolfi</i> . . . . .	281
4.1.4	Causality patterns of a marketing campaign conducted over time: evidence from the latent Markov model. <i>Fulvia Pennoni, Leo Paas and Francesco Bartolucci</i>	289

##### **4.2 Complex Spatio-temporal Processes and Functional Data** **297**

4.2.1	Clustering of spatio-temporal data based on marked variograms. <i>Antonio Balzanella and Rosanna Verde</i> . . . . .	297
4.2.2	Space-time earthquake clustering: nearest-neighbor and stochastic declustering methods in comparison. <i>Elisa Varini, Antonella Peresan, Renata Rotondi, and Stefania Gentili</i> . . . . .	304
4.2.3	Advanced spatio-temporal point processes for the Sicily seismicity analysis. <i>Marianna Silino and Giada Adelfio</i> . . . . .	312
4.2.4	Spatial analysis of the Italian seismic network and seismicity. <i>Antonino D’Alessandro, Marianna Silino, Luca Greco and Giada Adelfio</i> . . . . .	320

##### **4.3 Dimensional Reduction Techniques for Big Data Analysis** **328**

4.3.1	Clustering Data Streams via Functional Data Analysis: a Comparison between Hierarchical Clustering and K-means Approaches. <i>Fabrizio Mauro, Francesca Fortuna, and Tonio Di Battista</i> . . . . .	328
4.3.2	Co-clustering algorithms for histogram data. <i>Francisco de A.T. De Carvalho and Antonio Balzanella and Antonio Irpino and Rosanna Verde</i> . . . . .	338
4.3.3	A network approach to dimensionality reduction in Text Mining. <i>Michelangelo Misuraca, Germana Scepi and Maria Spano</i> . . . . .	344
4.3.4	Self Organizing Maps for distributional data. <i>Rosanna Verde and Antonio Irpino</i>	352

##### **4.4 Enviromental Processes, Human Activities and their Interactions** **353**

4.4.1	Estimation of coral growth parameters via Bayesian hierarchical non-linear models. <i>Crescenza Calculli, Barbara Cafarelli and Daniela Cocchi</i> . . . . .	353
4.4.2	A Hierarchical Bayesian Spatio-Temporal Model to Estimate the Short-term Effects of Air Pollution on Human Health. <i>Fontanella Lara, Ippoliti Luigi and Valentini Pasquale</i>	361
4.4.3	A multilevel hidden Markov model for space-time cylindrical data. <i>Francesco Lagona and Monia Ranalli</i> . . . . .	367
4.4.4	Estimation of entropy measures for categorical variables with spatial correlation. <i>Linda Altieri, Giulia Rolì</i> . . . . .	373

##### **4.5 Innovations in Census and in Social Surveys** **381**

4.5.1	A micro-based approach to ensure consistency among administrative sources and to improve population statistics. <i>Gianni Corsetti, Sabrina Prati, Valeria Tomeo, Enrico Tucci</i> . . . . .	381
4.5.2	Demographic changes, research questions and data needs: issues about migrations. <i>Salvatore Strozza and Giuseppe Gabrielli</i> . . . . .	392

4.5.3	Towards more timely census statistics: the new Italian multiannual dissemination programme. <i>Simona Mastroluca and Mariangela Verrascina</i> . . . . .	400
<b>4.6</b>	<b>Living Conditions and Consumption Expenditure in Time of Crises</b>	<b>409</b>
4.6.1	Household consumption expenditure and material deprivation in Italy during last economic crises. <i>Ilaria Arigoni and Isabella Siciliani</i> . . . . .	409
<b>4.7</b>	<b>Network Data Analysis and Mining</b>	<b>418</b>
4.7.1	Support provided by elderly Italian people: a multilevel analysis. <i>Elvira Pelle, Giulia Rivellini and Susanna Zaccarini</i> . . . . .	418
4.7.2	Data mining and analysis of comorbidity networks from practitioner prescriptions. <i>Giancarlo Ragozini, Giuseppe Giordano, Sergio Pagano, Mario De Santis, Pierpaolo Cavallo</i> . . . . .	426
4.7.3	Overlapping mixture models for network data (manet) with covariates adjustment. <i>Saverio Ranciati and Giuliano Galimberti and Ernst C. Wit and Veronica Vinciotti</i>	434
<b>4.8</b>	<b>New Challenges in the Measurement of Economic Insecurity, Inequality and Poverty</b>	<b>440</b>
4.8.1	Social protection in mitigating economic insecurity. <i>Alessandra Coli</i> . . . . .	440
4.8.2	Changes in poverty concentration in U.S. urban areas. <i>Francesco Andreoli and Mauro Mussini</i> . . . . .	450
4.8.3	Evaluating sustainability through an input-stateoutput framework: the case of the Italian provinces. <i>Achille Lemmi, Laura Neri, Federico M. Pulselli</i> . . . . .	458
<b>4.9</b>	<b>New Methods and Models for Ordinal Data</b>	<b>466</b>
4.9.1	Weighted and unweighted distances based decision tree for ranking data. <i>Antonella Plaia, Simona Buscemi, Mariangela Sciandra</i> . . . . .	466
4.9.2	A dissimilarity-based splitting criterion for CUBREMOT. <i>Carmela Cappelli, Rosaria Simone and Francesca Di Iorio</i> . . . . .	474
4.9.3	Constrained Extended Plackett-Luce model for the analysis of preference rankings. <i>Cristina Mollica and Luca Tardella</i> . . . . .	480
4.9.4	A prototype for the analysis of time use in Italy. <i>Stefania Capecchi and Manuela Michellini</i> . . . . .	487
<b>4.10</b>	<b>New Perspectives in Supervised and Unsupervised Classification</b>	<b>493</b>
4.10.1	Robust Updating Classification Rule with applications in Food Authenticity Studies. <i>Andrea Cappozzo, Francesca Greselin and Thomas Brendan Murphy</i> . . . . .	493
4.10.2	A robust clustering procedure with unknown number of clusters. <i>Francesco Dotto and Alessio Farcomeni</i> . . . . .	500
4.10.3	Issues in joint dimension reduction and clustering methods. <i>Michel van de Velden, Alfonso Iodice D'Enza and Angelos Markos</i> . . . . .	508
<b>4.11</b>	<b>New Sources, Data Integration and Measurement Challenges for Estimates on Labour Market Dynamics</b>	<b>514</b>
4.11.1	The development of the Italian Labour register: principles, issues and perspectives. <i>C. Baldi, C. Ceccarelli, S. Gigante, S. Pacini</i> . . . . .	514
4.11.2	Digging into labour market dynamics: toward a reconciliation of stock and flows short term indicators. <i>F. Rapiti, C. Baldi, D. Ichim, F. Pintaldi, M. E. Pontecorvo, R. Rizzi</i> . . . . .	523
4.11.3	How effective are the regional policies in Europe? The role of European Funds. <i>Gennaro Punzo, Mariateresa Ciommi, and Gaetano Musella</i> . . . . .	531
4.11.4	Labour market condition in Italy during and after the financial crises: a segmented regression analysis approach of interrupted time series. <i>Lucio Masserini and Matilde Bini</i> . . . . .	539

# Evaluating sustainability through an input-state-output framework: the case of the Italian provinces

## *Valutazione delle sostenibilità attraverso un sistema input-state-output: analisi sulle province Italiane*

Achille Lemmi<sup>1</sup>, Laura Neri<sup>2</sup>, Federico M. Pulselli<sup>3</sup>

**Abstract** In line with the recommendation of monitoring local context, in this paper we propose to investigate regional (NUTS2) and provincial (NUTS3) economic systems in a schematic and usable way: three different indicators are used to take into account resource use (input), societal organization (state) and to quantify the outputs of the system (output). A fuzzy cluster analysis is applied to the input-state-output indicator framework, that, as a whole, represents the interconnection of the three aspects of sustainability, namely environmental, social and economic. This framework is a useful and comprehensive tool for investigating and monitoring local context economic systems.

**Abstract** *In linea con le direttive di monitorare sistemi come le economie nazionali in un contesto locale, questo lavoro, propone l'analisi dei sistemi economici regionali (NUTS2) e provinciali (NUTS3) in modo semplice schematico e fruibile: tre diversi indicatori sono utilizzati per tenere conto dell'uso delle risorse (input), organizzazione della società (stato) e per quantificare l'output del sistema (output). Un'analisi cluster sfocata viene applicata agli indicatori del framework input-state-output, che, nel suo complesso, rappresentano l'interconnessione dei tre aspetti della sostenibilità, in particolare ambientale, sociale ed economica.*

**Key words:** Sustainability, input-state-output, fuzzy cluster analysis

---

<sup>1</sup> Achille Lemmi, ASESU Tuscan Universities Research Centre "Camilo Dagum", lemmiachille@virgilio.it

<sup>2</sup> Laura Neri, Department of Economics & Statistics, University of Siena, laura.neri@unisi.it

<sup>3</sup> Federico M. Pulselli, Ecodynamics Group, Department of Earth, Environmental and Physical Sciences, University of Siena, federico.pulselli@unisi.it

## 1 Introduction

In recent years there has been an increasing interest in the measurement of collective phenomena at the local level. The EU Committee of the Regions (2014) strongly recommend local authorities to define their own “2020 vision”, based on a territorial dimension, overcoming the present top-down approach of country targets fitting all regions irrespectively. Italy is subdivided into 20 regions (NUTS2) representing the first-level of administrative divisions; the country is further subdivided into 107 provinces (NUTS3). Though progressive measures are trying to eliminate this intermediate administrative level, provinces still play an important role in planning, coordination and cooperation at local level in connection with municipalities and other local bodies. According to OECD Regional Well-being (2016, <http://www.oecd.org/cfe/regional-policy/how-life-country-facts-italy.pdf>) “Italy has the largest regional disparities among the OECD countries in safety, with the Aosta Valley ranking in the top 1% and Sicily in the bottom 10% of the OECD regions. Important regional differences are found also in jobs, environment, community, civic engagement, income and access to services”. In line with the recommendation of monitoring local context, in this paper we propose to investigate regional (NUTS2) and provincial (NUTS3) economic systems in a schematic and usable way: three indicators are used to take into account resource use, societal organization and to quantify the outputs of the system. This framework is consistent with an input-state-output scheme (I-S-O, Pulselli et al., 2015), representing the ordered triad environment–society–economy. A three-storey pyramid represents the mutual relationships among the three dimensions of sustainability, rotating the pyramid clockwise, the succession of the stages is oriented from left to right, consistently with the I-S-O framework

## 2 Data and Methods

In this framework different combinations of indicators can be used to account for the input-state-output. The study here presented is referred to provincial areas, so the preliminary challenge to face is the data availability. Then, the aim of the research is to produce an “objective” classification of the Italian provinces in terms of the three aspects of sustainability. Such classification should be useful for designing and delivering policy responses to economic, environmental and social needs at the local level.

### 2.1 *The indicators*

The input indicator should be representative of what a system extracts/obtains, directly or indirectly, from the environment. Referring to provincial areas, which are sub-regional systems, poor datasets are systematically produced, especially in the environmental field. Therefore, no encompassing methods, like energy evaluation (as in Pulselli et al., 2015) or ecological footprint, or other



environmental accounting methods can apply. In this case we used an aggregation of energy consumption measures, collected from two institutional databases (Terna: National electric network; DGERM: Ministry of Economic Development). In particular, we selected electricity consumption and sale of a set of fuel types, for all the provinces of Italy. In order to aggregate these measures, we calculated the equivalent in terms of CO<sub>2</sub> emission to show both the use of resources (electricity and fuels) and the environmental pressure (emissions) on the other. The result is an estimation of gross CO<sub>2</sub> emission due to almost all the items composing the energy sector. In order to monitor the environmental pressure of human activities on each provincial territory and compare provinces, the amount of CO<sub>2</sub> per unit area is computed. This choice helps determine the contribution of human actions (that imply energy consumption) to climate change independently of the number of inhabitants in each area. To encompass the characteristics of the state, a synthetic indicator describing a form of societal organization should be used. Considering the critical importance of reducing unemployment, in order to drive toward inclusive society, an indicator related to the labour market seems to be appropriate. Again, availability and reliability of data at NUTS3 level is the critical issue to face: it has been retained that the most reliable official statistics should be the Labour Force Survey, so the unemployment rate has been chosen as key state indicator. Gross Domestic Product (GDP) is the principal aggregate for measuring economic development/growth of a country/region. In this analysis we maintain this logic, considering the GDP per inhabitants in purchasing parity power as output indicator of the I-S-O system.

## 2.2 *Methods*

Our goal is to produce an “objective” classification of the Italian provinces in terms of the three aspects of sustainability, namely environmental, social and economic, according to the chosen triad: CO<sub>2</sub> per unit area, unemployment rate and GDP per inhabitants in purchasing parity power.

Data clustering is recognized as a statistical technique for classifying data elements into different groups (known as clusters) in such a way that the elements within a group possess high similarity while they differ from the elements in a different group. By using the triad of indicators, the clustering algorithm starts from an initial partition of the provinces into a fixed number of groups, where each group is initially randomly chosen. The grouping is then updated: based on the distance between every single observation and the reference objects of each group, every observation is reallocated to the closest group aiming to produce a classification that is reasonably “objective” and “stable”. The classification obtained by using the crisp cluster analysis suffers for both a poor homogeneity within group and a lacking separation between the groups. For this reason, we explored the clustering procedure by using a soft clustering known as Fuzzy Cluster Analysis, a very important clustering technique based on fuzzy logic. In case of soft clustering techniques, fuzzy sets are used to cluster data, so that each point may belong to two or more clusters with different degrees of membership. In many situations, fuzzy clustering is more natural than hard clustering.

Objects on the boundaries between several clusters are not forced to fully belong to one of the cluster, but rather are assigned membership degrees between 0 and 1 indicating their partial membership. The most popular algorithm, the fuzzy c-means (developed by Dunn in 1973 and improved by Bezdek in 1981), aims at minimizing an objective function –the weighted within-groups sum of square- whose (main) parameters are the membership degrees and the parameters determining the localisation as well as the shape of the clusters. Objects are assigned to clusters according to membership degrees in  $[0,1]$ : 0 is where the data point is at the farthest possible point from a cluster's center and 1 is where the data point is the closest to the center. Each of the  $c$  clusters is represented by a cluster center. These centers are chosen randomly in the beginning, then each data vector is assigned to the nearest prototype according to a suitable similarity measure and each center is replaced by the centre of gravity of those data assigned to it. The alternating assignment of data to the nearest center and the update of the cluster centres is repeated until the algorithm converges, i.e., no more changes happen.

Although the extension from crisp to fuzzy clustering seems to be an obvious concept, it turns out that to actually obtain membership degrees between zero and one, it is necessary to introduce a so-called fuzzifier in fuzzy clustering. Usually, the fuzzifier is simply used to control how much clusters are allowed to overlap. This fuzzifier function creates an area of crisp membership values around a prototype while outside of these areas of crisp membership values, fuzzy values are assigned. The analysis has been performed by using the fuzzy clustering with polynomial fuzzifier (Frank Klawonn and Frank Hoppner, 2003).

A problem that frequently occurs in real data analysis is the presence of one or more observation presenting anomalous values, i.e. outliers. Such a subset, that may be referred to as noise, tends to disrupt clustering algorithms making difficult to detect the cluster structure of the remaining domain points. According to the adopted approach the first  $k$  standard clusters are homogeneous, whereas the noise cluster, serving as a “garbage collector”, contains the outliers and is usually not formed by objects with homogeneous.

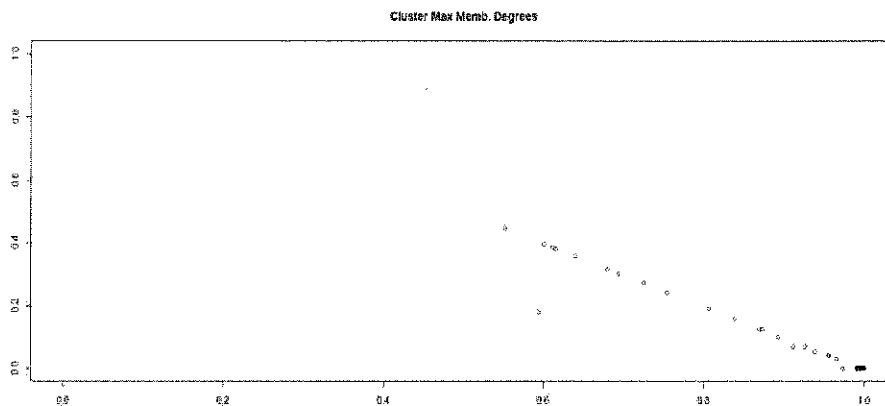
All data have been standardized before performing the cluster analysis. The analysis has been conducted by using R and specifically the R package *fclust* (Giordani, Ferraro, 2018). The package provides the cluster solution, cluster validity index and plots and also the visualization of fuzzy clustering results.

The analysis has been performed by using the fuzzy clustering with polynomial fuzzifier with noise cluster. For assessing cluster validity, some have been evaluated: here, just the so called partition coefficient (PC) is presented. Given that the closer to unity the PC index the “crisper” the clustering is and that value close to  $1/n_c$  (where  $n_c$  is the number of clusters) indicates that there is no clustering tendency, a  $PC=0.92$ , indicates a significant cluster structure.

As regard to the visual inspection of fuzzy clustering results -VIFCR- (Klawonn et al., 2003) is a scatter plot where, for each object, the coordinates  $u_1$  and  $u_2$  denote, respectively, the highest and the second highest membership degrees. All points lie within the triangle with vertices  $(0,0)$ ,  $(0.5,0.5)$  and  $(1,0)$ . In the ideal case of (almost) crisp membership degrees all points are near the vertex  $(1,0)$ . This graph has been

evaluated for different partition and the partition with three cluster, plus the noise one, containing just one province (Milan), seems to be the best one (Fig.1).

The algorithm applied to reach the fuzzy clustering solution, assigns an objects to clusters only if the corresponding member function degree is greater than 0.5. In this way the closest hardest partition can be identified by assigning an object to the cluster according to the maximal membership function ( $>0.5$ ) and the characteristics of each cluster can be identified. The cardinalities and the average membership function of the closest hardest cluster are reported in Table 1, as well as the average of each indicator considered in the clustering procedure. Specifically, in Table 2, the list of provinces belonging to each cluster, according to the maximal membership function, are reported; it is worth to point out that cluster 4, the noise one, includes just Milan.



**Figure 1:** Scatter plot: for each observation, the coordinates  $u_1$  and  $u_2$  denote, respectively, the highest and the second highest membership degrees.

**Table 1:** Size, average membership function and average values for the indicators by cluster

<i>Cluster</i>	<i>size</i>	<i>m.f.</i>	<i>CO<sub>2</sub> area</i>	<i>Unempl rate</i>	<i>pps ab</i>
1	27	0.96	2272.88	7.75	30866.67
2	44	0.94	912.13	10.31	24854.34
3	35	0.97	804.13	19.78	17046.71
4	1	1	12755.85	7.68	44493.13

In order to have a vision of the distribution of each indicator within each cluster, the boxplots can be observed (Figure:2).

PCA plot (Figure 3) is a very useful tools to visualize the data: this has nothing to do with the type of clustering algorithm or the accuracy of the algorithm used, however it is a useful representation to recognize the utility of the fuzzy clustering, given that the clusters, are clearly separated but the borderline units provinces in this study, are very close.

**Table 2:** The list of provinces belonging to the four clusters according to the maximal membership function

<i>Cluster</i>	<i>Provinces</i>
1	Aosta,Bergamo,Bologna,Bolzano,Brescia,Como,Cremona,Firenze, Forli-Cesena, Genova,Lecco,Livorno,Mantova, Modena,Padova, Parma,Prato,Ravenna,Reggio-Emilia,Roma,Trento,Treviso,Trieste,Varese,Venezia,erona,Vicenza
2	Alessandria,Ancona,Arezzo,Ascoli,Piceno,Asti,Belluno,Biella,Chieti,Cuneo, Ferrara,Frosinone,Gorizia,Grosseto,Imperia,L'Aquila,La Spezia, Latina,Lodi, Lucca,Macerata,Massa Carrara,Novara,Nuoro,Pavia,Perugia,Pesaro Urbino, Pescara,Piacenza,Pisa,Pistoia,Pordenone,Potenza,Rieti,Rimini,Rovigo,Savona, Siena,Sondrio,Teramo,Terni,TorinoUdine,Verbano-Cusio-Ossola,Vercelli
3	Agrigento,Avellino,Bari,Benevento,Brindisi,Cagliari,Caltanissetta,Campobasso, Carbonia-L.,Caserta,Catania,Catanzaro,Cosenza, Crotone,Enna,Foggia,Isernia, Lecce,Matera,Medio,Campidano,Messina,Napoli,Ogliastra,Olbia-T.,Oristano, Palermo,Ragusa,Reggio C., Salerno,Sassari,Siracusa,Taranto,Trapani,Vibo Valentia,Viterbo
4	Milano

### 3. An attempt of taxonomy according to the triad

Cluster composition (Table 2) and average values of indicators (Tab. 1) suggest a high heterogeneity among clusters emphasizing the existence of economic disparities. In Cluster 1 there are 27 provinces: Rome, plus 26 provinces located in the North or Central Italy. Comparing Cluster 1 with Cluster 2 and 3, we can observe the highest average level of CO<sub>2</sub> per area, the highest GDP per capita and the lowest unemployment rate. Observing Fig. 2, we can also realize that the unemployment rate is quite homogeneous within the cluster while the GDP per capita and the level of CO<sub>2</sub> per area present the largest variability if compared with the variability of Cluster 2 and 3. This result is justified considering the presence of very peculiar provinces in this cluster, like Bolzano (belonging to an autonomous region) with the highest GDP per capita (40000 Euro) (out of Milan), and the lowest unemployment rate among all the Italian provinces, and the lowest level of CO<sub>2</sub> per area within Cluster 1. and Livorno, the province presenting the lower membership function for this cluster (0.64), maybe due to its GDP which is closer to the average GDP per capita computed for Cluster 2. Cluster 2 is the most heterogeneous cluster as regard to the geography of provinces. It is composed mainly by provinces located in the North and Central Italy plus provinces of Lazio (out of Rome and Viterbo) and Abruzzo plus Nuoro and Potenza. All the provinces in Cluster 2, out of Torino, have small medium dimension. This cluster is characterized by high variability in GDP, ranging from 17900 for Nuoro to 29900 for Siena and low variability in CO<sub>2</sub>. Torino presents the lowest membership function (0.55) among all the provinces, being a borderline unit between Cluster 2 and 1. Cluster 3 is composed by whole Calabria, Campania, Molise, Puglia, Sicilia and Sardegna (out of Nuoro), plus Viterbo (Lazio). This cluster is very similar to the one, labelled as “Minimal system with high social risks” by Bertin (2012) in his classification concerning the welfare state systems of the Italian regions. Cluster 3 is characterized by the lowest average level of CO<sub>2</sub> emission per area, although such level is biased

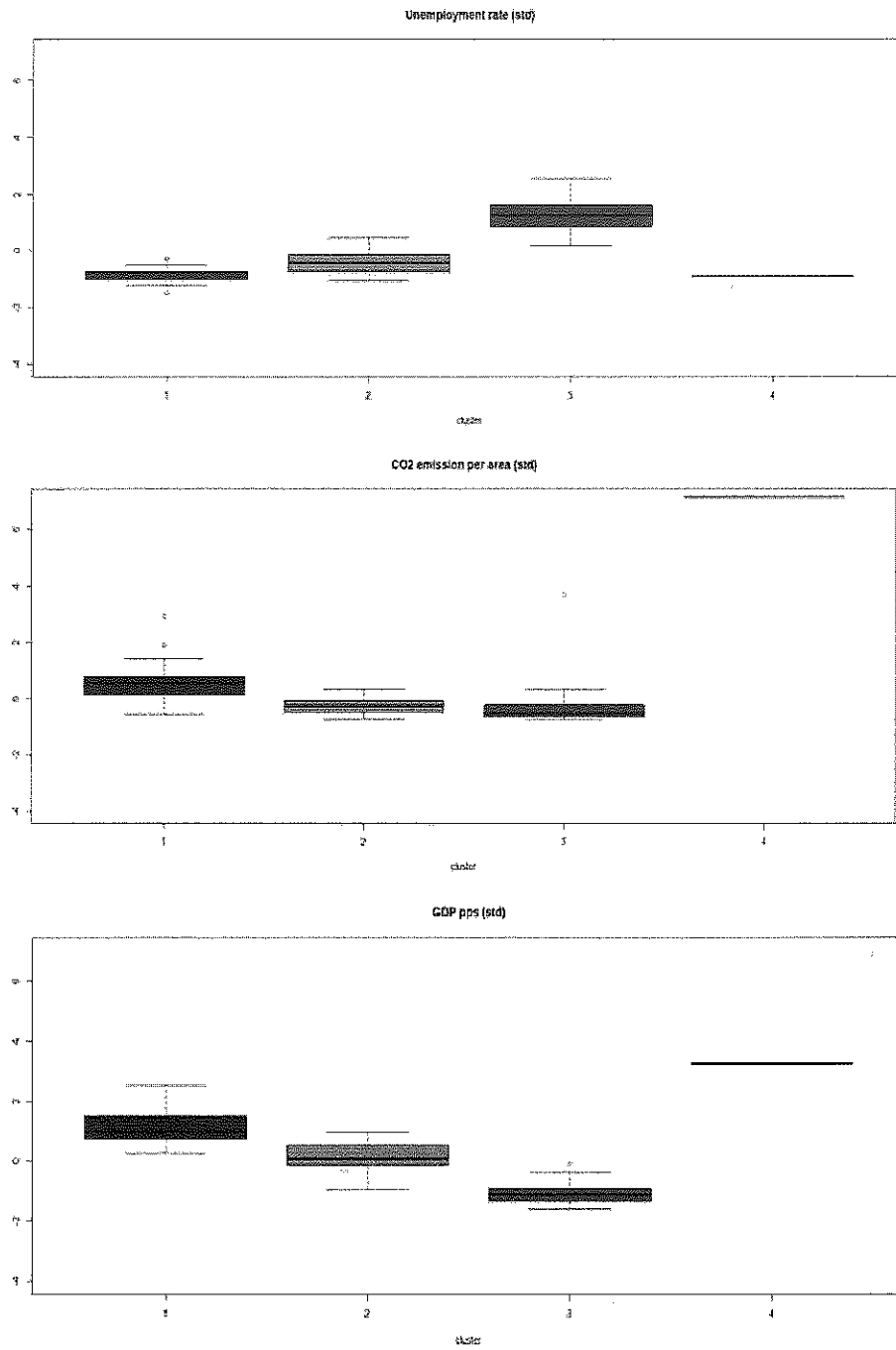


Figure 2: Boxplot-cluster structure

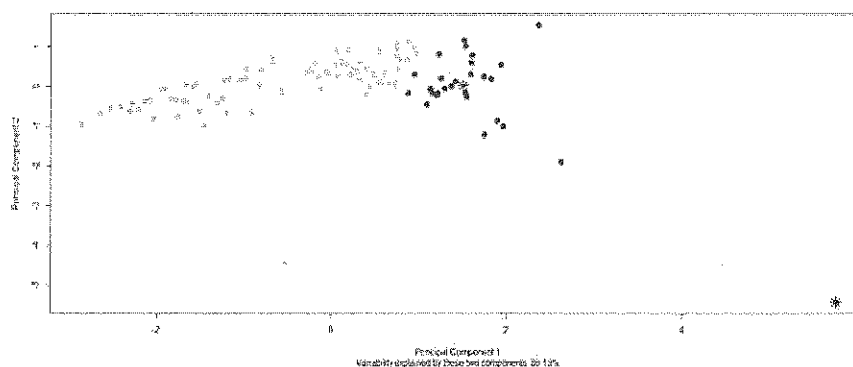


Figure 3: Principal component plot of the cluster structure

due to the membership of Napoli, presenting the highest level of CO<sub>2</sub> among all Italian provinces (out of Milan). As regard to the unemployment rate and GDP, the cluster presents an evident disparity with respect to the others. As regard to GDP per capita, the level of Cagliari (23400), the highest value within cluster 3, is lower than the minimum value in Cluster 1 (25500) registered for Livorno. It is interesting to observe that Napoli is a borderline unit between cluster 3 and 2.

Finally, Cluster 4, the one containing just Milan, is for sure a true outlier with respect to the level of CO<sub>2</sub> emission per area (more than ten times the average values of all the other provinces i.e. 12755.85 vs 1218.29) and GDP per capita which in Milan is nearly double than the average values of all the other provinces (44493 vs 23812).

The analysis conducted confirms the well-known dualism, resulted in a North-South divide in GDP per capita and in labour-market performance, adding a new element: the Southern Italian provinces are homogeneous with respect to the considered characteristics whilst the Northern and Central provinces are not homogeneous even if they belong to the same region. This result suggests that local policies can be better aligned and tailored to specific local opportunities and challenges. Moreover, the clustering structure obtained can be an useful tool to adopt, considering that provinces belonging to the same cluster could share ways to generate better outcomes for environment, jobs and the economy, trying to reach provinces of their cluster.

## References

1. Bertin, G.: Crises and change processes of the welfare systems, *Italian Sociological Review*, vol. 2, n. 1, pp. 1–13 (2012)
2. Bezdek, J.C.: *Pattern Recognition with Fuzzy Objective Function Algorithms*, Kluwer Ac. Pub., Norwell, MA, (1981)
3. Dunn, J.C.: A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters, *J. Cybern.* 3, 3, 32-57 (1974)
4. Klawonn, F., Höppner, F.: What is fuzzy about fuzzy clustering? Understanding and improving the concept of the fuzzifier. In: *Advances in intelligent data analysis*, pp. 254-264. LNCS (2003)
5. Giordani, P., Ferraro, M.B.: *fclust: Fuzzy Clustering*, R package version 1.1.2, <http://cran.r-project.org/package=fclust>.
6. Pulselli, F.M., Coscieme, L., Neri, L., Regoli, A., Sutton, P., Lemmi, A., Bastianoni, S.: The World Economy in a cube: a more rational structural representation of Sustainability. *Glob. Environ. Change* 35, pp. 41–51 (2015).