

# Robust multivariate analysis for mixed-type data: novel algorithm and its practical application in socio-economic research

Aurea Grané<sup>a</sup>, Silvia Salini<sup>b,\*</sup>, Elena Verdolini<sup>c</sup>

<sup>a</sup>*Universidad Carlos III de Madrid*

<sup>b</sup>*University of Milan*

<sup>c</sup>*University of Brescia and RFF-CMCC European Institute of Economics and Environment, Centro Euro-Mediterraneo sui Cambiamenti Climatici*

---

## Abstract

We propose a novel method and algorithm for the analysis and clustering of mixed-type data using a hierarchical approach based on Forward Search. In our procedure, the identification of groups is based on the identification of similar trajectories and then linked to very intuitive two-dimensional maps. The proposed algorithm can use different measures for the calculation of distance in the case of mixed-type data, such as Gower's metric and Related metric scaling. A key feature of our algorithm is its ability to discard redundant information from a given set of variables. The practical usefulness of the algorithm is illustrated through two applications of high relevance for empirical economic research. The first one focuses on comparing different indicators of environmental policy stringency in different countries. The second one applies our procedure to identify clusters of countries based on information regarding their institutional characteristics.

*Keywords:* Forward Search, Mixed Type Data, Outliers, Robustness, Redundant information, Clustering

---

---

\*Corresponding author, Department of Economics, Management and Quantitative Methods, Via Conservatorio 7, 20122 Milano, Italy, Tel: +390250321538

*Email addresses:* agrane@est-econ.uc3m.es (Aurea Grané), silvia.salini@unimi.it (Silvia Salini), elena.verdolini@unibs.it (Elena Verdolini)

## 1. Introduction

Mixed-type data comprise both numeric and categorical features, and mixed datasets frequently occur in many domains, such as economics, health, finance, marketing, including data coming from socio-demographic surveys. Applied economists and social scientists are often faced with the necessity to deal with mixed-type data. For instance, mixed data indicators measuring a given economic or societal aspect often need to be compared to understand the extent to which they convey similar or different information, as in Galeotti et al. (2020). Furthermore, clustering often needs to be applied to mixed datasets to find structures and to group similar observations for further analysis, as in Nesta et al. (2019). These contributions highlight the challenges associated with the use of mixed-type data for socio-economic research. To begin with, one cannot rely on a simple distance measure, such as the Euclidean distance, because of the presence of categorical data. Moreover, in the statistical literature a few distance measures to deal with mixed data exist, such as Gower's similarity index (Ahmad and Khan, 2019) but they are plagued by important shortcomings, as highlighted in Grané and Romera (2018) and discussed more at length below. See also Foss et al. (2019) for clustering methods for mixed data and van de Velden et al. (2018) for distance based methods for mixed data.

A recent relevant methodological contribution in the context of mixed data is presented in Grané and Romera (2018), who construct robust profiles for mixed-type data using multidimensional scaling, which is one of the most extended methodologies to visualize the profile structure of mixed data. To this end, Grané and Romera (2018) compare different multidimensional scaling configurations (coming from different metrics) through a combination of sensitivity and robust analysis. They propose a robust joint metric combining different distance matrices, avoiding redundant information, via Related Metric Scaling (RelMS) as an alternative to classical Gower's metric.

The first (methodological) contribution of this paper to the literature on mixed data is the development of a novel robust algorithm for the explanatory data analysis of mixed datasets. This is achieved by combining the related metric scaling measure proposed by Grané and Romera (2018) with a Forward Search algorithm (Atkinson and

31 Riani (2004)). On the one hand, related metric scaling allows to overcome the main  
32 shortcomings of Gower’s measure. On the other hand, Forward Search (FS) is a pow-  
33 erful general method which can be applied to many statistical models to make them  
34 robust. The FS algorithm was introduced by Atkinson and Riani (2000), (2004) in the  
35 context of robust regression models and has been extended to many other fields, such  
36 as financial models, cluster analysis, curve monitoring, robust inference, and such. In  
37 our context, Forward Search is useful because (a) it incorporates flexible data-driven  
38 trimming for the detection of outliers and unsuspected structure in the data and (b) it  
39 facilitates data visualization, in particular it allow us to visually represent how the pro-  
40 cedure to calculate the related metric scaling joint metric unfolds rather than providing  
41 only a final picture of the outcome.

42 The second (practical) contribution of this paper is to demonstrate the usefulness of  
43 this novel algorithm for applied socio-economic analysis through two empirical appli-  
44 cations of applied economic analysis. First, we show the usefulness of our approach for  
45 the comparison of mixed-data indicators of environmental policy which underline the  
46 analysis of Galeotti et al. (2020). We demonstrate the need of an alternative measure  
47 to Gower’s metric in presence of mixed-type data by showing how RelMS can discard  
48 redundant information from different indicators. Furthermore, we use a stability anal-  
49 ysis to show how the Multi Dimensional Scaling (MDS) configurations of RelMS are  
50 more stable than those using Gower.

51 Second, we apply our method to the widely known dataset described in La Porta et  
52 al. (1999) to show its usefulness in generating clusters for countries in terms of the key  
53 institutional dimensions. This procedure can be used to generate an index of similarity  
54 for potential use in applied research, such as the generation of instrumental variables  
55 similarly to what proposed in Nesta et al. (2019). The La Porta et al. (1999) database  
56 is a mixed dataset containing variables describing four key country-level institutional  
57 aspects: legal origin, political freedom, efficiency of institutions and interference with  
58 the private sector. These are important underlying characteristics of a country’s insti-  
59 tutional, legal and political framework.

60 The rest of the paper is organized as follows: Section 2 is devoted to describe the  
61 proposed algorithm. In Section 3 we present an alternative metric more robust than

62 Gower's. Section 4 presents the application to the two empirical applications related  
63 to environmental policy stringency and to institutional aspects: we describe the data,  
64 apply our algorithm, and comment the results and their usefulness for applied socio-  
65 economic research. Section 5 concludes, highlighting other potential application areas  
66 and discussing future research avenues.

## 67 **2. Method**

68 In this Section, we first describe the Forward Search Distance Based (FS-DB) al-  
69 gorithm. This novel approach combines the FS method with a distance-based tool,  
70 used in Grané and Romera (2018) to detect outliers in mixed-type datasets. While this  
71 distance-based tool can cope with any distance measure, the algorithm is initially de-  
72 scribed in terms of Gower's distance, since we are interested in mixed-type datasets.  
73 A more interesting alternative is given in Section 3, where the distance is tailored via  
74 RelMS.

### 75 *2.1. The Foward Search philosophy of data analysis.*

76 The FS is a data-driven strategy which is based on carefully chosen subsets of  
77 the data. The key difference with respect to other robust strategies for data analy-  
78 sis, is that the algorithm is not only based on one subsample, but on a sequence of  
79 subsets of the original data. It is an adaptive hard trimming method (Salini et al.,  
80 2016). In the words of their initial proponents "the FS is not a simple new algo-  
81 rithm but a new philosophy of looking at the data, which involves watching a film  
82 of the data rather than a snapshot". The crucial idea behind the FS approach is to  
83 monitor how a fitted model changes whenever a new statistical unit is added to the  
84 subset. The model of interest is initially fitted on a starting subset, whose units can  
85 change in each step of the algorithm. Thus, this approach helps to understand the ef-  
86 fect that each unit (outlier or not, leverage point or not) exerts on the fitted model (see  
87 <http://rosa.unipr.it/FSDA/guide.html> for a more detailed description  
88 of the method).

89 2.2. The FS-DB algorithm

90 The idea behind our proposed approach is to help understand the structure of mixed-  
 91 type datasets by identifying the subset of closest units (according to a user-selected  
 92 distance measure) as well as those units that are the most distant from the set(s) of the  
 93 data. Apart from the numerical outputs, there are two graphical outputs of our algo-  
 94 rithm. First, the Forward-plot with the trajectories of the units which illustrate their  
 95 performance along the steps of the algorithm. Second, the MDS-plot with the final  
 96 MDS representations of the dataset. Graphical outputs are explained in Section 4. The  
 97 algorithm implemented code has been submitted in order to be included in the next re-  
 98 lease of the common and flexible framework provided by the FSDA Toolbox of Matlab  
 99 (Riani et al. (2012)<sup>1</sup>)

100 The starting point of our procedure is a data matrix of mixed type of dimension  
 101  $n \times p$ . The steps we follows are:

- 102 1. Select a distance measure. In this first example we use Gower’s similarity coef-  
 103 ficient. Given two  $p$ -dimensional vectors  $\mathbf{z}_i$  and  $\mathbf{z}_j$ , Gower’s similarity coefficient  
 104 is defined as

$$s_{ij} = \frac{\sum_{h=1}^{p_1} (1 - |z_{ih} - z_{jh}|/R_h) + a + \alpha}{p_1 + (p_2 - d) + p_3}, \quad 0 \leq s_{ij} \leq 1, \quad (1)$$

105 where  $p = p_1 + p_2 + p_3$ ,  $p_1$  is the number of continuous (or quantitative) vari-  
 106 ables,  $a$  and  $d$  are the number of positive and negative matches, respectively, for  
 107 the  $p_2$  binary variables,  $\alpha$  is the number of matches for the  $p_3$  multi-state cate-  
 108 gorical variables, and  $R_h$  is the range of the  $h$ -th continuous variable. Gower’s  
 109 distance is defined as  $\delta^2(\mathbf{z}_i, \mathbf{z}_j) = 1 - s_{ij}$ , which are the entries of the matrix of  
 110 squared distances  $\Delta$ .

- 111 2. Select a subset size ( $m < n$ ). By default  $m$  is set as 10% of  $n$ .  
 112 3. Select the units inside the starting subset which have lowest distance measure.  
 4. Calculate the geometric variability of the subset  $V_{\Delta(m)}$ . Let  $\{\mathbf{z}_i, 1 \leq i \leq m\}$  be  
 $m$   $p$ -dimensional vectors containing the information of the  $m$  individuals in the

---

<sup>1</sup>The FSDA Toolbox of Matlab is freely available at <http://rosa.unipr.it/fsddownload.html>.

subset and consider a matrix  $\Delta_{(m)}$  of squared distances, with entries  $\delta^2(\mathbf{z}_i, \mathbf{z}_j)$ , for  $1 \leq i, j \leq m$ . The geometric variability of  $\Delta_{(m)}$  is

$$V_{\Delta_{(m)}} = \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m \delta^2(\mathbf{z}_i, \mathbf{z}_j).$$

5. Calculate for each unit outside the subset the *distance-based proximity* function  $\phi(i)$  to the subset. Given a new individual  $\mathbf{z}_0 \in \mathbb{R}^p$ , the distance-based proximity of  $\mathbf{z}_0$  to the set  $\{\mathbf{z}_i, 1 \leq i \leq m\}$  is

$$\phi(\mathbf{z}_0) = \frac{1}{m} \sum_{i=1}^m \delta^2(\mathbf{z}_0, \mathbf{z}_i) - V_{\Delta_{(m)}}.$$

- 113 6. Include in the subset the unit with the minimum value of  $\phi(i)$ ; set  $m$  equal to  
 114  $m + 1$ .  
 115 7. Iterate the procedure from step 3 until all  $n$  units are included in the subset.  
 116 8. Monitoring  $\phi(i)$  for each unit on the subset size.  
 117 9. Plot the trajectory in multidimensional scaling (MDS) maps and identify groups  
 118 and outliers.

119 In this implementation we select the units inside the starting subset which lowest  
 120 distance measure. However, note that step 3 allows the units to enter and exit the subset,  
 121 since in each iteration the current subset is formed by those units with lowest distance  
 122 measure. Another interesting approach, that we can explore for future development,  
 123 is that detailed in Atkinson et al. (2006) where units in the initial subset are randomly  
 124 chosen in order to check the stability to the starting point.

### 125 3. An alternative to Gower's metric: Related metric scaling

126 Gower's similarity coefficient is one of the most popular similarity measures and  
 127 perhaps the easiest way to obtain a distance measure when working with mixed-type  
 128 data. However, it presents two important drawbacks. The first one, pointed out long  
 129 time ago by Gower (1992); Krzanowski (1994), is that, just like any distance function  
 130 satisfying additivity with respect to variables, this coefficient ignores any association

131 (correlation) between variables and, thus, is not able to discard any redundant infor-  
 132 mation. The second drawback is the lack of robustness: this coefficient uses the stan-  
 133 dardized city block distance for quantitative variables (see equation (1)), which is not a  
 134 robust measure. As a consequence, in the presence of outliers, the stability of the MDS  
 135 configurations can be affected, as shown in Grané and Romera (2018). This second  
 136 drawback may be solved by replacing standardized city block distance by, for instance,  
 137 a robustified Mahalanobis distance. However, still the first drawback will remain in the  
 138 new coefficient. Thus, our proposal is to overcome both shortcomings by obtaining a  
 139 distance measure for mixed-type data via related metric scaling.

140 Related metric scaling (RelMS) is a multivariate technique that allows to obtain a  
 141 unique representation of a set of observations from several distance matrices computed  
 142 on the same set of observations. The method is based on the construction of a joint  
 143 metric that satisfies several axioms related to the property of identifying and discarding  
 144 redundant information (Cuadras (1998); Cuadras and Fortiana (1998)).

145 Given a set of  $k \geq 2$  matrices of squared distances measured on the same group of  $n$   
 146 observations,  $\{\Delta_\alpha\}_{\alpha=1,\dots,k}$ , the first requirement in the construction of the joint metric  
 147 is that all matrices  $\Delta_\alpha$  have the same geometric variability. Note that the condition of  
 148 equal geometric variability can always be assumed to hold, since multiplying a squared  
 149 distances matrix by an appropriate constant amounts to a change of measurement unit.

150 • First step: Standardize each matrix  $\Delta_\alpha$  with respect to its geometric variability  
 151  $V_{\Delta_\alpha} = \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n \delta^2(\mathbf{z}_i, \mathbf{z}_j)$ , where  $\delta^2(\mathbf{z}_i, \mathbf{z}_j)$  are the entries of matrix  $\Delta$ , for  
 152  $1 \leq i, j \leq n$ . In an abuse of notation, we call  $\Delta_\alpha$  its standardized version.

• Second step: For each distance matrix  $\Delta_\alpha$  consider its doubly centered inner  
 product matrix:

$$\mathbf{G}_\alpha = -\frac{1}{2} \mathbf{H} \Delta_\alpha \mathbf{H},$$

153 where  $\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}\mathbf{1}'$ ,  $\mathbf{I}_n$  is the identity matrix of order  $n$  and  $\mathbf{1}$  is a  $n \times 1$  vector  
 154 of ones.

• Third step: Compute the inner product matrix of the joint metric as

$$\mathbf{G} = \sum_{\alpha=1}^k \mathbf{G}_\alpha - \frac{1}{k} \sum_{\alpha \neq \beta} \mathbf{G}_\alpha^{1/2} \mathbf{G}_\beta^{1/2}, \quad (2)$$

155 where  $\mathbf{G}_\alpha^{1/2}$  denotes the square root of  $\mathbf{G}_\alpha$ , which can be obtained through the  
156 singular value decomposition of  $\mathbf{G}_\alpha$ .

- Fourth step: The matrix of the joint metric can be computed as

$$\Delta = \mathbf{g}\mathbf{1}' + \mathbf{1}\mathbf{g}' - 2\mathbf{G}, \quad (3)$$

157 where  $\mathbf{g} = \text{diag}(\mathbf{G})$  is a column vector containing the diagonal of matrix  $\mathbf{G}$ .

158 Note that, in order to obtain MDS configurations, it is enough to work with formula  
159 (2), although formula (3) is required for computing  $\phi(i)$  in the Forward Search.

160 How do we interpret formula (2)? The first addend of this formula mimics Gower's  
161 similarity coefficient by adding the  $k$  metrics; the second one is responsible of discard-  
162 ing redundant information coming from different sources. This second addend provides  
163 more flexibility to the MDS configuration when working with mixed-type data. Thus,  
164 RelMS allows us to tailor a metric to reflect specific information of a mixed dataset.  
165 See Grané and Romera (2018) for the mathematical properties of the method. Another  
166 advantage of the method is that it does not require data preprocessing.

167 In our application, we construct the joint metric from  $k = 3$  distance matrices, one  
168 for each variable type. In particular,  $\Delta_1$  contains the information related to quantitative  
169 variables and we use a robust version of Mahalanobis distance (Riani et al., 2009) to  
170 compute it. For multi-state categorical variables, we start by computing the similarity  
171 matrix  $\mathbf{S}_2$  with Sokal-Michener's pairwise similarity coefficient (matching coefficient),  
172 and then we obtain  $\Delta_2 = 2(\mathbf{1}\mathbf{1}' - \mathbf{S}_2)$ ; Finally, for binary variables, we compute the  
173 similarity matrix  $\mathbf{S}_3$  with Jaccard's pairwise similarity coefficient and we get  $\Delta_3 =$   
174  $2(\mathbf{1}\mathbf{1}' - \mathbf{S}_3)$ .

#### 175 **4. Application**

176 The lack of a robust method for mixed-data is clearly an important limitation for  
177 applied research in environmental and sustainability issues, as argued above. In this  
178 section, we specifically illustrate this point through two examples. First, the lack of  
179 a robust approach to deal with mixed data is evident in the analysis of Galeotti et al.



180 (2020), who explore and compare different proxies of environmental policy stringency  
181 in a sample of 189 OECD countries over the years 1995-2009. Their analysis shows  
182 that different indexes of environmental policy stringency can give rise to significantly  
183 different country rankings (and, implicitly, country clustering) depending on whether  
184 they are based on data of environmental expenditures (i.e. inputs), emissions (i.e. out-  
185 puts) or of the type of policy instrument implemented. The extent to which different  
186 indicators provide similar (and thus redundant information) is a matter of concern. Fur-  
187 thermore, the inability to identify outliers in a multivariate framework implies that any  
188 index or summary statistics (in this case, of environmental policy stringency) will not  
189 be robust, rather it will be influenced by such outliers and possibly mask their presence.  
190 This is due to the fact that composite indicators are necessarily data driven. Therefore,  
191 a major implication of the analysis presented in Galeotti et al. (2020) is the lack of a  
192 robust method to detect outliers and to analyze differences and similarities in a mixed  
193 data context.

194       Second, the recent contribution of Nesta et al. (2019) highlights the importance  
195 of being able to cluster countries based on several characteristics, some of which can  
196 be easily measured on continuous scales while others are necessarily summarized by  
197 categorical or binary proxies. Specifically, Nesta et al. (2019) use information on in-  
198 stitutional policies to generate an Instrumental Variable in the context of their main  
199 research question, namely the impact of environmental policy on the direction of in-  
200 novation. Specifically, Nesta et al. (2019) use information on the underlying legal and  
201 institutional characteristics of their sample of countries in order to cluster them. For  
202 each country, they then use the data on environmental policy stringency of all other  
203 observations in the cluster to generate the instrument used to address the endogeneity  
204 of their main variable of interest. A robust way to provide clear country clustering  
205 in this context would be really useful to address the issue of endogeneity in a more  
206 satisfactory way.

207       Using these two specific cases as illustrative, in the rest of this Section we show  
208 the main advantages of our novel algorithm, and its potential for improving the use of  
209 mixed-type data in socio-economic analysis.

210 *4.1. Mixed data on environmental policy stringency: comparing different indicators*

211 The aim of this subsection is to demonstrate both the need for and the usefulness  
212 of our approach for the comparison of mixed-data indicators of environmental policy  
213 stringency. To do so, we use the mixed data which underline the analysis of Galeotti  
214 et al. (2020), namely different indicators of environmental policy stringency for 18  
215 OECD countries in the year 2009.<sup>2</sup> Descriptive statistics of the variables used are  
216 provided in Table 1, alongside the original data source. Figure 1 provides an overview  
217 of the values of the different indicators in our sample. As highlighted in Galeotti et  
218 al. (2020), the different indicators sometimes provide similar information regarding  
219 the level of environmental policy stringency in a given country (for instance, in the  
220 case of *NOx* Taxes and *SOx* Taxes for most of the countries), while in other cases the  
221 information provided is very different (for instance, this is the case with *FITs* – Feed  
222 in Tariffs for solar and wind on the one hand and with *CO<sub>2</sub>* Tradable Certificates).

223 The first step to demonstrate the need for and usefulness of our approach as alter-  
224 native to Gower’s for mixed-type data is to compare the MDS configurations computed  
225 from Gower’s and RelMS metrics in this context. As Figure 2 shows, the percentage of  
226 explained variability is greater when using RelMS metric. In both configurations there  
227 are four countries quite far from the others (Canada, USA, Austria and Turkey). How-  
228 ever, the relative positions of the other countries are different in the two configurations.  
229 For example, with Gower’s metric Japan looks close to Denmark, whereas when using  
230 RelMS Japan is more isolated.

231 *4.1.1. Percentage of redundant information*

232 One of the attractive properties of RelMS is the ability to discard redundant infor-  
233 mation. For this reason, we calculate the percentage of redundant information in the

---

<sup>2</sup>Note that this application is done on a sub-set of the original Galeotti et al. (2020) data, namely a cross-section of the original database which contained data for the 19 OECD countries over the years 1995–2009. The countries included in this analysis are Austria, Australia, Canada, Germany, Denmark, Spain, Finland, France, the United Kingdom, Greece, Hungary, Italy, Japan, Netherlands, Portugal, Sweden, Turkey, and the United States

Variable	Name	Mean	Std. Dev.	Min.	Max.	Type	Data Source
CO2 Tax	<i>CO2 – Tax</i>	0.333	1.414	0	6	categorical, ordered	Botta and Kozluz (2014)
NOx Tax	<i>NOx – Tax</i>	1.667	2.169	0	6	categorical, ordered	Botta and Kozluz (2014)
SOx Tax Indicator	<i>SOx – Tax</i>	1.444	1.977	0	6	categorical, ordered	Botta and Kozluz (2014)
CO2 Certificates	<i>CO2 – TraS</i>	3.222	2.157	0	6	categorical, ordered	Botta and Kozluz (2014)
Green Certificates	<i>Green – TraS</i>	1.167	1.855	0	6	categorical, ordered	Botta and Kozluz (2014)
White Certificates	<i>White – TraS</i>	0.611	1.29	0	5	categorical, ordered	Botta and Kozluz (2014)
FIT Wind Indicator	<i>Wind – FIT</i>	1.833	1.581	0	5	categorical, ordered	Botta and Kozluz (2014)
FIT Solar Indicator	<i>Solar – FIT</i>	2.944	2.235	0	6	categorical, ordered	Botta and Kozluz (2014)
Sulphur Content Indicator	<i>Sulph – cont</i>	5.833	0.383	5	6	categorical, ordered	Botta and Kozluz (2014)
R&D Indicator	<i>RD – indicator</i>	2.722	1.904	0	6	categorical, ordered	Botta and Kozluz (2014)
Diesel Tax	<i>Diesel – tax</i>	4.111	1.132	2	6	categorical, ordered	Botta and Kozluz (2014)
DRS Indicator	<i>DRS – indicator</i>	0.556	0.511	0	1	binary	Botta and Kozluz (2014)
NOx Limits	<i>NOx – Limits</i>	4.333	1.749	1	6	categorical, ordered	Botta and Kozluz (2014)
SOx Limits	<i>SOx – Limits</i>	4.389	1.614	0	6	categorical, ordered	Botta and Kozluz (2014)
PM Limits	<i>PM – Limits</i>	2.222	1.263	1	6	categorical, ordered	Botta and Kozluz (2014)
Levinson Indicator EM	<i>BL – EM</i>	1.151	0.306	0.685	1.719	continuous	Galeotti et al. (2019)
Levinson Indicator CO2	<i>BL – CO2</i>	1.23	0.454	0.721	2.245	continuous	Galeotti et al. (2019)
Levinson Indicator SOX	<i>BL – SOX</i>	2.492	2.745	0.211	11.958	continuous	Galeotti et al. (2019)
Levinson Indicator NOX	<i>BL – NOX</i>	1.43	1.318	0.197	6.263	continuous	Galeotti et al. (2019)
Levinson Indicator NMVOC	<i>BL – NMVOC</i>	1.748	1.172	0.369	4.466	continuous	Galeotti et al. (2019)
Levinson Indicator NH3	<i>BL – NH3</i>	1.288	0.875	0.574	4.302	continuous	Galeotti et al. (2019)
Levinson Indicator N2O	<i>BL – N2O</i>	1.43	1.098	0.482	4.25	continuous	Galeotti et al. (2019)
Levinson Indicator CO	<i>BL – CO</i>	2.635	1.541	0.118	5.885	continuous	Galeotti et al. (2019)
Levinson Indicator CH4	<i>BL – CH4</i>	1.699	1.352	0.581	6.761	continuous	Galeotti et al. (2019)
Energy R&D Intensity	<i>RDD – GDP</i>	0.519	0.48	0.029	1.898	continuous	IEA (2015)

Table 1: Variable description and descriptive statistics for Environmental Stringency data. *Categorical and binary variables come from Botta and Kozluz (2014). Continuous variables come from Galeotti et al. (2019) and are computed following the approach proposed in Brunel and Levinson (2013) using data on emissions and value added from WIOD 2013 (Timmer 2015). The last variable is computed using data on Energy R&D Investments from the IEA (2015) and data on GDP from the National Accounts of OECD Countries (both in constant 2013 PPP US dollars)*

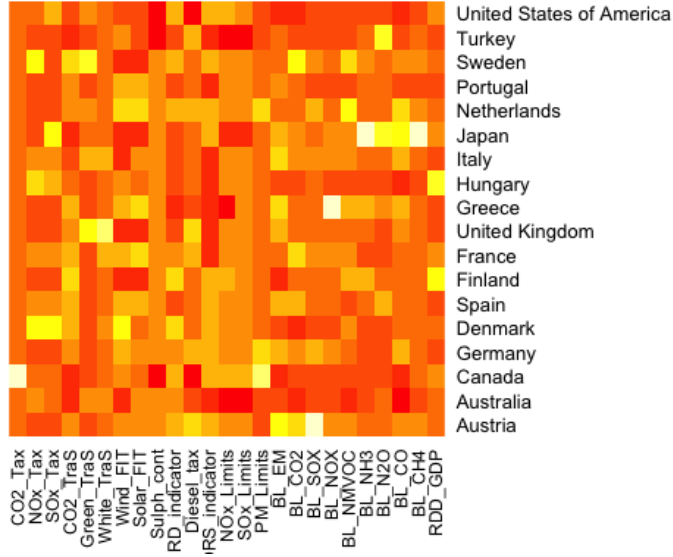


Figure 1: Heatmap of Environmental Stringency data

234 indicators for the countries in our sample. This will also help to explain why the MDS  
 235 configurations obtained with Gower’s metric or RelMS are rather different.

We first of all need to understand whether matrices  $\Delta_1$ ,  $\Delta_2$ , and  $\Delta_3$  share a high percentage of information. To answer this question we compute the following measure of agreement between two matrices:

$$\rho(\alpha, \beta) = \frac{V(\Delta_\alpha, \Delta_\beta)}{\|\Delta_\alpha\| \|\Delta_\beta\|},$$

236 where  $V^2(\Delta_\alpha, \Delta_\beta) = \frac{1}{n} |\text{tr}(\mathbf{G}_\alpha) - \text{tr}(\mathbf{G}_\beta)|$ ,  $\|\Delta_\alpha\|^2 = \text{tr}(\mathbf{G}_\alpha)$ . Coefficient  $\rho$  takes values  
 237 in  $[0, 1]$ , being equal to one in case of orthogonality (that is, the Euclidean configura-  
 238 tions associated to  $\Delta_\alpha$  and  $\Delta_\beta$  generate orthogonal subspaces on  $\mathbb{R}^n$ ) and equal to zero  
 239 in case of equality (that is, in case  $\Delta_\alpha = \Delta_\beta$ ). Then, the percentage of information  
 240 shared by two distance matrices is obtained as  $(1 - \rho)100$  percent.

241 In our case, the percentages of shared information by matrices  $\Delta_1$ ,  $\Delta_2$ , and  $\Delta_3$  are  
 242 shown in Table 2, where we can see that these matrices contain redundant information  
 243 (indeed they share more than 85% of the information). This is one of the reasons why  
 244 the MDS configurations look different when using RelMS metric or Gower’s. A second

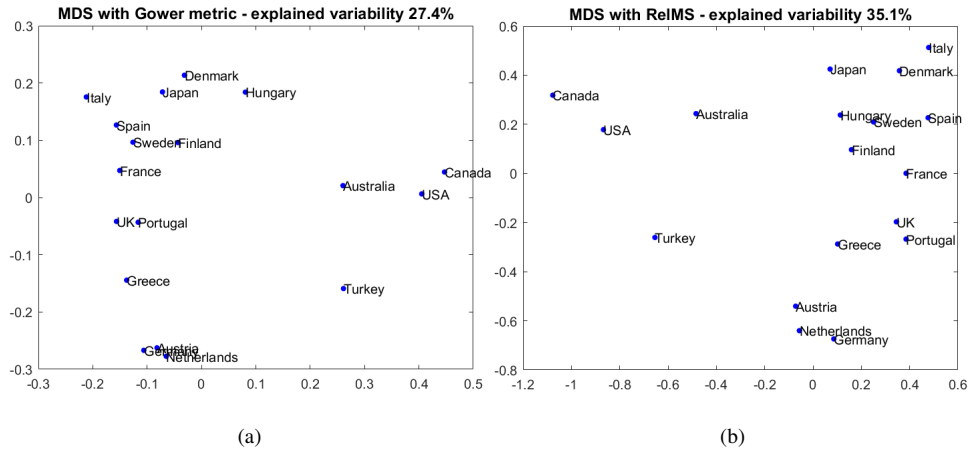


Figure 2: MDS configurations obtained from (a) Gower's metric and (b) RelMS metric

245 reason explaining why the configurations look different is related with the stability of  
 246 those configurations, as we discuss next.

Table 2: Percentages of information shared by matrices  $\Delta_1$ ,  $\Delta_2$ , and  $\Delta_3$

	$\Delta_1$	$\Delta_2$	$\Delta_3$
$\Delta_1$	100%	90.0%	90.4%
$\Delta_2$		100%	86.2%
$\Delta_3$			100%

247 *4.1.2. Stability analysis*

248 Here we are interested in evaluating the influence of the  $i$ th observation on the other  
 249  $n - 1$  observations, in the sense that how the exclusion of the  $i$ th observation from the  
 250 original dataset can affect the MDS configuration of the  $n - 1$  remaining points. To  
 251 solve this question we apply the leave-one-out cross-validation procedure proposed by  
 252 Krzanowski (2006).

253 The idea of this method is to leave out each observation (that is, each row of the  
 254 dataset) in turn and to compute the MDS configuration with the remaining  $n - 1$  obser-

255 vations. Then, the excluded observation is projected onto the MDS configuration using  
256 Gower's interpolation formula (Gower, 1968) leading to an "augmented" configuration.  
257 Finally, the  $n$  "augmented" configurations are compared with the original one (that is,  
258 that obtained from the whole dataset) by superimposing them. More specifically, they  
259 are just put on top each other (correctly aligned), as described in Krzanowski (2006),  
260 with no Procrustean rotation. Since sometimes the  $n(n + 1)$  points may overload the  
261 diagram, it is recommended to surround each point with the smallest hypersphere that  
262 contains a given percentage (e.g., 95 percent) of the cross-validatory replicate points.  
263 Hence, small hyperspheres indicate a very stable point, whereas large ones a very un-  
264 stable one.

265 In Figure 3 we depict the 95 percent-stability regions for the MDS configurations  
266 using Gower's metric (a-panel) and RelMS metric (b-panel). The radius of each hyper-  
267 sphere is given by the squared root of the 95th quantile of the  $l^2$  distances between the  
268 original coordinates of the point and the coordinates of its replicates. We can see that in  
269 panel (b) there is only one point very unstable (influential), whereas in panel (a) most  
270 of the points are unstable (and thus, influential). Additionally, we can compute some  
271 descriptive statistics for the radii. For example, the mean (and SD) are 0.1617 (0.1683)  
272 for Gower's metric and 0.0096 (0.0684) for RelMS metric. Hence, our conclusion is  
273 that MDS configuration computed from RelMS distance is more stable than Gower's.

#### 274 4.1.3. *Forward Search trajectory of the two metrics*

275 Figure 4 below shows the Forward plot using Gower's index (a-panel) and RelMS  
276 (b-panel). In this figure, trajectories which end close to one another represent countries  
277 which are similar among themselves, but different from others.

278 As expected the two measures produce different trajectories: we know from the  
279 previous analysis that Gower's metric does not discard redundant information. The  
280 brushing units, highlighted in red, are the units that enter at the end of the search, so  
281 the most distant from the bulk of the data. In this case, it is apparent that relying on  
282 Gower's metric suggests that Australia, Canada, Japan, Turkey and United States of  
283 America may cluster together, and separately from other countries (a panel). When  
284 accounting for redundant information, on the other hand, these countries do not appear

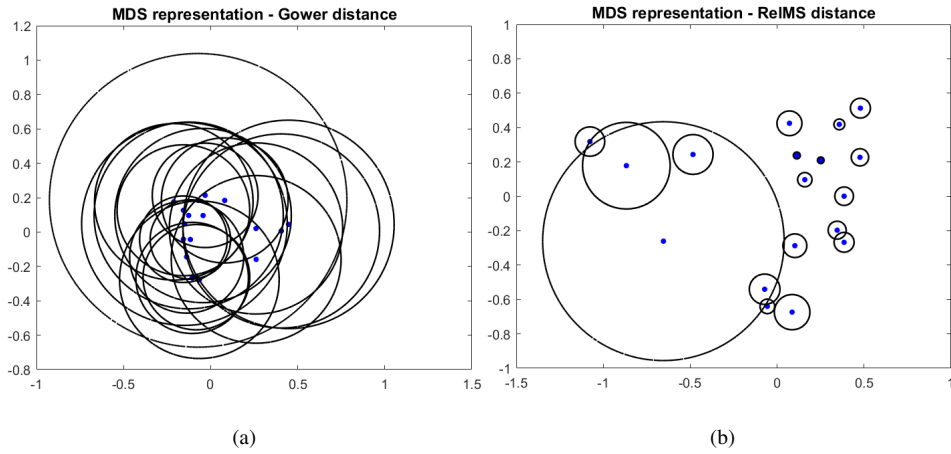


Figure 3: Sensitivity analysis of MDS configurations. In (a) sensitivity analysis of the MDS configuration from Gower’s metric; in (b) sensitivity analysis of the MDS configuration from RelMS metric

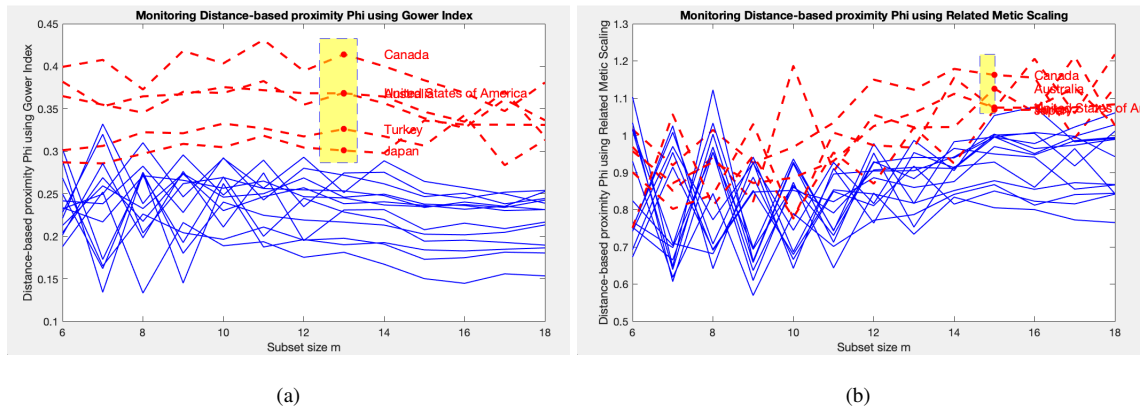


Figure 4: Forward plots: monitoring the  $\phi$  distance based proximity measure. In (a) you can see the trajectories using Gower’s measure, in (b) you can see the trajectories using Related Metric Scaling

285 as different from the others (b panel). We further discuss the usefulness of our approach  
 286 for clustering purposed in the next subsection, which contains our second application.

287 *4.2. Mixed data on countries’ institutional structure: clustering*

288 In order to demonstrate the usefulness of our novel algorithm for clustering obser-  
 289 vations, we rely on data on institutional aspects which we source our data from the

290 widely known dataset described in La Porta et al. (1999). This mixed dataset contains  
291 several variables describing four key country-level institutional aspects: legal origin,  
292 political freedom, efficiency of institutions and interference with the private sector.  
293 The “Legal origin” variables identify the legal origin of the Company Law or Com-  
294 mercial law of a given country. They are a set of binary indicators identifying if a  
295 country is of either British, French, Socialist, German and Scandinavian legal origin.  
296 “Political Freedom” is measured with two proxies: a democracy index and a politi-  
297 cal freedom index, both ranging from 0 to 10. Lower values indicate lower levels of  
298 political freedom. “Efficiency of Institutions” is measured through three variables: cor-  
299 ruption, bureaucratic delays and tax compliance. Corruption and bureaucratic delays  
300 range from 0 to 10, while tax compliance is measured on the scale from 0 to 6. In all  
301 three cases, the index increase when efficiency increases. Lastly, interference with the  
302 private sector is measured with an index of property rights and a business regulation  
303 index, both ranging from 1 to 5. For a more detailed description of the database, please  
304 refer to La Porta et al. (1999). For tractability, we limit our dataset to a sample of 35  
305 countries.<sup>3</sup>

306 All these variables selected for our analysis were chosen because they can help to  
307 identify economies with similar underlying institutional structures. Identifying clusters  
308 of countries is indeed potentially relevant for the study of economic outcomes, on the  
309 one hand, and for use to generate instrumental variables in economic analyses. Indeed,  
310 each of these variables separately has been proposed as, or used in the computation  
311 of, instrumental variable in previous literature (as, for instance, in Nesta et al. (2019)).  
312 Yet, the lack of an aggregation technique appropriate for mixed data meant that each  
313 variable could only be used separately, and that the usefulness of focusing on different  
314 institutional aspects has remained largely unexplored.

315 Figure 5 below shows the Forward plot using Gower index (a-panel) and using

---

<sup>3</sup>The countries included in our sample are: Australia, Austria, Belgium, Brazil, Canada, Switzerland, China, Czech Republic, Germany, Denmark, Spain, Finland, France, United Kingdom, Greece, Hungary, Indonesia, India, Ireland, Iceland, Italy, Japan, Republic of Korea, Mexico, Netherlands, Norway, Poland, Portugal, Russian Federation, Slovak Republic, Slovenia, Sweden, Turkey, United States, South Africa.



316 related metric scaling (b-panel). In this figure, trajectories which end close to one  
 317 another represent countries which are similar among themselves, but different from  
 318 others. The two measures produce different trajectories. A main reason for this is the  
 319 fact that Gower's metric does not discard redundant information and it is a less robust  
 320 measure with respect to related metric scaling. That is, individuals that look close with  
 321 Gower's metric, may not look so close when using a robust metric that can discard  
 322 redundant information.

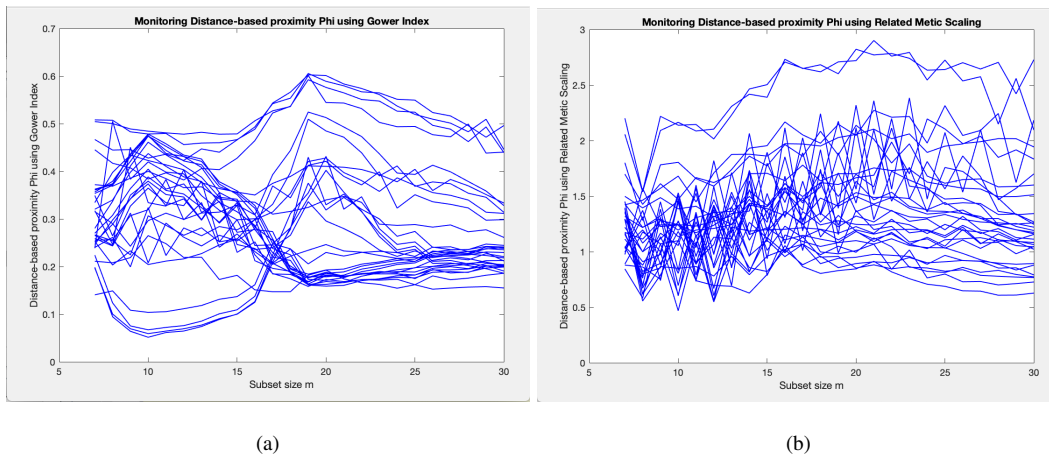


Figure 5: Forward plots: monitoring the  $\phi$  distance based proximity measure. In (a) you can see the trajectories using Gower's measure, in (b) you can see the trajectories using Related Metric Scaling

323 In Figure 6 and Figure 7 we show some groups of countries with similar trajectories  
 324 in the related metric scaling configuration, the more robust one. Using the implemented  
 325 algorithms it is possible to brush the trajectory and to show the selected countries in the  
 326 multidimensional scaling maps. In this example we consider three coordinates but it is  
 327 possible to plot the units in a different space. For example in the scatter-plot matrix of  
 328 the original quantitative variable, coloring of splitting for level of categorical or binary  
 329 variables.

330 Figure 6 highlights the units that enter at the begin of the search, so the units which  
 331 are nearest to one another in terms of the different indicators of environmental pol-  
 332 icy stringency. These countries are Austria, Canada, Spain, Greece, Ireland, Norway,

333 Portugal. Figure 7 highlights the units that enter at the end of the search, so the most  
 334 distant from the bulk of the data, China, Indonesia, India and Mexico. This evidence  
 335 confirm that these fast developing countries are very close when it comes to their en-  
 336 vironmental policy stringency. Importantly, note that the order in which observations  
 337 enter the search are not determined by the level of the mixed data considered, rather by  
 338 the distance of difference observations from one another.

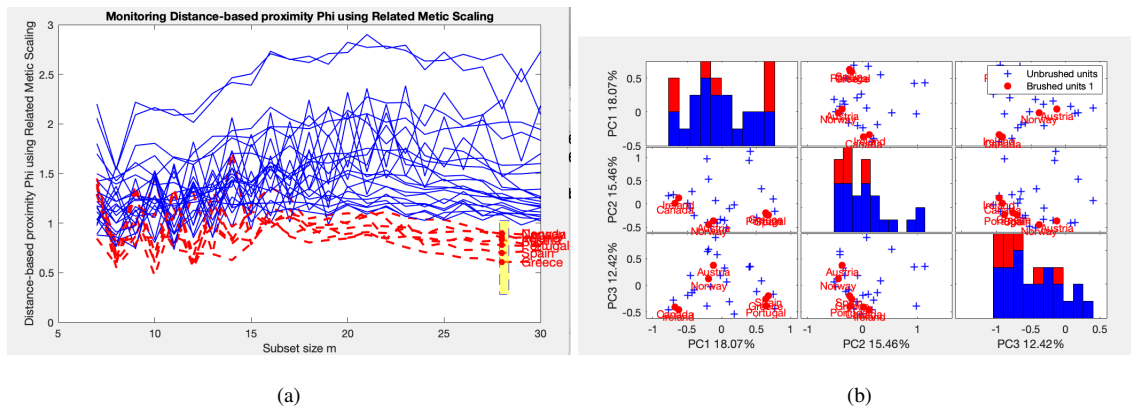


Figure 6: Output of the algorithm using the joint metric obtained via related metric scaling. In (a) you can see Forward plot, in (b) you can see the MDS Plot

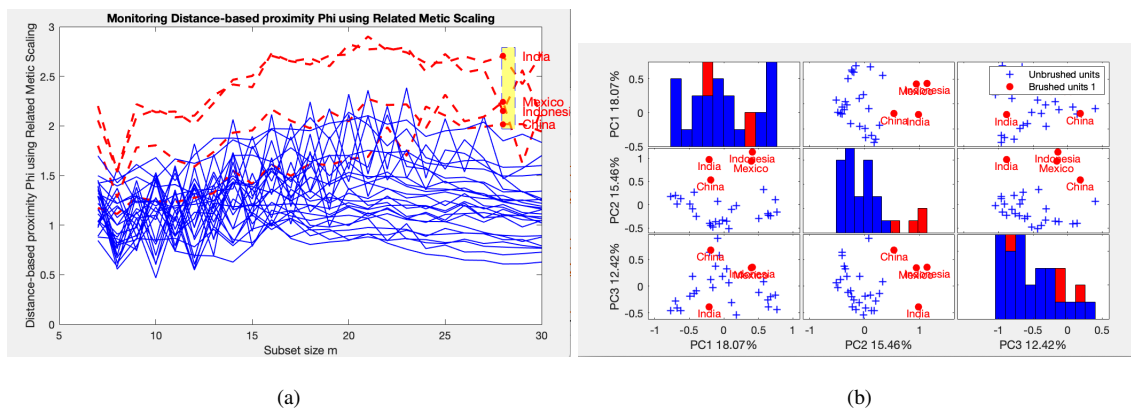


Figure 7: Output of the algorithm using the Related Metric Scaling. In (a) you can see Forward plot with highlighted some units that enter at the end of the search using the brushing option, in (b) you can see the corresponding MDS Plot

339 This algorithm and the procedure to handle the clustering of various observations  
340 over step in the forward search clearly offer the practical advantages of providing a  
341 clustering when in presence of mixed data. It also speaks to the potential of using an  
342 underlying score of similarity against these broad-based aspects of institutional design  
343 and quality. This could potentially be of great use for applied socio-economic re-  
344 searchers to generate exogenous instrumental variables using information about other  
345 countries in the cluster to instrument for one’s own variable of interest, following the  
346 procedure of Nesta et al. (2019)

## 347 **5. Conclusions**

348 This paper develops a novel robust algorithm for the explanatory data analysis of  
349 mixed datasets. We code this approach in the common and flexible computational  
350 framework provided by the FSDA Toolbox of Matlab by combining the RelMS joint  
351 metric proposed by Grané and Romera (2018) with a Forward Search algorithm.<sup>4</sup>  
352 From the methodological point of view, this is a significant improvement for two  
353 reasons. On the one hand, the related metric scaling allows to overcome the main  
354 shortcomings of Gower’s measure, which up to recently has been the most common  
355 approach to the analysis of mixed datasets. On the other hand, applying the Forward  
356 Search method we incorporate in our procedure flexible data-driven trimming for the  
357 detection of outliers and unsuspected structure in the data and we facilitate data visu-  
358 alization. Another advantage is that the method does not requires data preprocessing.

359 Our analysis also points to fruitful avenues of future methodological research.  
360 These include, as just mentioned, the possibility to select the units inside the start-  
361 ing subset in the Forward Search randomly in order to check the stability to the starting  
362 point, as suggerer in Atkinson et al. (2006). Moreover new interactive options for data  
363 visualization to improve the brushing of the units, for instance to produce the scatter  
364 plot matrix of the original quantitative variables or to color the dots differently based  
365 on nature of the variables selected; the implementation of other robust distances, such

---

<sup>4</sup>The code is available under request to the authors and is in the process of optimization and checking, with the aim of adding it to the next FSDA release

366 as a robust Gower measure (Mahalanobis instead of Manhattan) in the second step of  
367 the algorithm; the optimization of the proposed methods for a larger datasets, includ-  
368 ing tests for anomalous data and contaminations. We are also exploring the possibility  
369 to develop the classical hierarchical cluster using Releted Metric Scaling as a distance  
370 measure inside the common and flexible computational framework provided by the  
371 FSDA Toolbox of Matlab, including the generation of a dendrogram. Finally, in sec-  
372 tion 4.1.2 we use the leave-one-out cross-validation procedure proposed by Krzanowski  
373 (2006) to check the stability of MDS. An interesting future development, also to be im-  
374 plemented in the FSDA toolbox, could be to apply the FS for the same purpose, with  
375 the aim to avoid the typical masking effect of the outliers.

376 The usefulness of this new method is illustrated through two applications relevant  
377 for applied socio-economic analysis. First, we build on Galeotti et al. (2020) and use  
378 our method to compare different indicators of environmental policy stringency. Sec-  
379 ond, we apply our novel approach to the widely known dataset of La Porta et al. (1999):  
380 we use data on institutional characteristics of a given country to generate country clus-  
381 ters which account for several complementary aspects, namely legal origin, political  
382 freedom, efficiency of institutions and interference with the private sector. These ex-  
383 amples confirm the high potential applicability of our novel approach beyond current  
384 applications for the clustering of observations and the generation of similarity indexes,  
385 including the generation of instrumental variables.

## 386 **6. Acknowledgments**

387 Elena Verdolini gratefully acknowledges funding from the Horizon 2020 Research  
388 and Innovation Programme under grant agreement No 730403 (INNOPATHS).

## 389 **References**

- 390 Ahmad A, Khan S (2019) Survey of State-of-the-Art Mixed Data Clustering Algo-  
391 rithms. *IEEE Access* 7: 31883-31902
- 392 Atkinson AC and Riani M (2000) *Robust Diagnosis Regression Analysis*. New York:  
393 Springer.

- 394 Atkinson A and Riani M (2004) The forward search and data visualisation, *Computational Statistics* 19: 29–54  
395
- 396 Atkinson A, Riani M, Cerioli A (2006) Random start forward searches with envelopes  
397 for detecting clusters in multivariate data. In *Data analysis, classification and the*  
398 *forward search*. Springer, Berlin, Heidelberg. 163–171)
- 399 Atkinson AC, Riani M, Cerioli A (2010) The forward search: Theory and data analysis.  
400 *Journal of the Korean Statistical Society* 39(2):117–134
- 401 Botta, E, Koźluk, T (2014) Measuring environmental policy stringency in OECD coun-  
402 tries: a composite index approach. OECD Economics Department Working Paper N.  
403 1177, <http://dx.doi.org/10.1787/5jxrjnc45gvg-en>
- 404 Brunel C, Levinson A (2013) Measuring environmental regulatory stringency. OECD  
405 Working Paper N. 2013/05. <http://dx.doi:10.1787/18166881>.
- 406 Cuadras CM (1998) Multidimensional dependencies in classification and ordination.  
407 *Analyses Multidimensionnelles des Données* pp 15–25
- 408 Cuadras CM, Fortiana J (1998) Visualizing categorical data with related metric scaling.  
409 In: *Visualization of Categorical Data*, Elsevier, pp 365–376
- 410 Cuadras CM, Fortiana J, Oliva F (1997) The proximity of an individual to a population  
411 with applications in discriminant analysis. *Journal of Classification* 14(1):117–136
- 412 Foss HA, Markatou M, Bonnie, R (2019) Distance Metrics and Clustering Methods for  
413 Mixed-type Data. *International Statistical Review* 87(1): 80–109.
- 414 Galeotti M, Salini S, Verdolini E (2020) Measuring Environmental Policy Stringency:  
415 Approaches, Validity, and Impact on Energy Efficiency. *Energy Policy* 136:111052.
- 416 Gower JC (1992) Generalized Biplots. *Biometrika* 79(4)75–93
- 417 Gower JC (1968) Adding a Point to Vector Diagrams in Multivariate Analysis.  
418 *Biometrika* 55:582–85.

- 419 Grané A, Romera R (2018) On visualizing mixed-type data: A joint metric approach  
420 to profile construction and outlier detection. *Sociological Methods & Research*  
421 47(2):207–239
- 422 IEA - International Energy Agency (2015). *Energy Technology RD&D Statistics*  
423 Database, available at [www.oecd.org](http://www.oecd.org).
- 424 Krzanowski WJ (1994) Ordination in the Presence of Group Structure for General Mul-  
425 tivariate Data. *Journal of Classification* 11:195–297
- 426 Krzanowski WJ (2006) Sensitivity in Metric Scaling and Analysis of Distance. *Bio-*  
427 *metrics* 62:239–44.
- 428 La Porta R, Lopez-de Silanes F, Shleifer A, Vishny R (1999) The quality of govern-  
429 ment. *The Journal of Law, Economics, and Organization* 15(1):222–279
- 430 Nesta, L, Verdolini E, Vona F (2019) Threshold policy effects and directed technical  
431 change in Energy Innovation. *Updated version of Nesta et al. 2018, Sciences Po*  
432 *publications No 2018–05*.
- 433 Riani M, Perrotta D, Torti F (2012) FSDA: a matlab toolbox for robust analysis  
434 and interactive data exploration. *Chemometrics and Intelligent Laboratory Systems*  
435 116:17–32
- 436 Riani M, Atkinson A, Cerioli A (2009) Finding an unknown number of multivariate  
437 outliers. *Journal of the Royal Statistical Society: series B (statistical methodology)*  
438 71.2:447–466
- 439 Salini S, Cerioli A, Laurini F, Riani M (2016) Reliable robust regression diagnostics.  
440 *International Statistical Review* 84(1):99–127
- 441 Timmer MP, Dietzenbacher E, Los B, Stehrer R, de Vries G.J. (2015) An Illustrated  
442 User Guide to the World Input–Output Database: the Case of Global Automotive  
443 Production. *Review of International Economics* 23:575–605
- 444 van de Velden M, Iodice D’Enza A and Markos A (2018) Distance-based clustering of  
445 mixed data. *Wires Computational Statistics*, 11(3).