

# Scaled Radial Axes for Interactive Visual Feature Selection: A Case Study for Analyzing Chronic Conditions

A. Sanchez<sup>a,d</sup>, C. Soguero-Ruiz<sup>a</sup>, I. Mora-Jiménez<sup>a</sup>, F. J. Rivas-Flores<sup>b</sup>, D. J. Lehmann<sup>c</sup>,  
M. Rubio-Sánchez<sup>a</sup>

<sup>a</sup>Universidad Rey Juan Carlos (Madrid, Spain)

<sup>b</sup>Hospital Universitario de Fuenlabrada (Madrid, Spain)

<sup>c</sup>Otto von Guericke University Magdeburg (Magdeburg, Germany)

<sup>d</sup>Research Center for Computational Simulation (Madrid, Spain)

---

## Abstract

In statistics, machine learning, and related fields, feature selection is the process of choosing a smaller subset of features to work with. This is an important topic since selecting a subset of features can help analysts to interpret models and data, and to decrease computational runtimes. While many techniques are purely automatic, the data visualization community has produced a number of interactive approaches where users can make decisions taking into account their domain knowledge. In this paper we propose a new visualization technique based on radial axes that allows analysts to perform feature selection effectively, in contrast to previous radial axes methods. This is achieved by employing alternative scaled axes that provide insight regarding the features that have a smaller contribution to the visualizations. Therefore, analysts can use the technique to carry out interactive backwards feature elimination, by discarding the least relevant features according to the information on the plots and their expertise. Our approach can be coupled with any linear dimensionality reduction method, and can be used when performing analyses of cluster structure, correlations, class separability, etc. Specifically, in this paper we focus on combining the proposed technique with methods designed for classification. Lastly, we illustrate the effectiveness of our proposal through a case study analyzing high-dimensional medical chronic conditions data. In particular, clinicians have used the technique for determining the most important features that discriminate between patients with diabetes and high blood pressure.

---

*Email addresses:* alberto.sanchez@urjc.es (A. Sanchez), cristina.soguero@urjc.es (C. Soguero-Ruiz), inmaculada.mora@urjc.es (I. Mora-Jiménez), franciscojavier.rivas@salud.madrid.org (F. J. Rivas-Flores), dirk@isg.cs.uni-magdeburg.de (D. J. Lehmann), manuel.rubio@urjc.es (M. Rubio-Sánchez)

*Keywords:* High-dimensional data visualization, interactive feature selection, visual analytics, exploratory data analysis, medical chronic conditions.

---

## 1. Introduction

The analysis of high-dimensional data sets is a complex and common problem in fields such as statistics, data mining, or machine learning. In practice, data sets may contain hundreds or thousands of features, many of which can be irrelevant, redundant, or simply add noise. Feature selection consists of the process of discarding those features. The topic is important since analyzing or using the resulting smaller subset can provide several benefits such as: simpler models that are easier to interpret, reduced overfitting, enhanced performance, or shorter computational runtimes.

While many feature selection techniques rely on purely automatic procedures (Guyon & Elisseeff, 2003), the data visualization community has produced a number of interactive approaches where users are integrated into the analysis process with the goal of benefiting from their perceptual capabilities, flexibility, and domain knowledge. With these visualization tools analysts are able to steer the selection process according to their expertise, obtaining subsets of features adapted to the specific problem and application domain, in contrast to automatic methods.

In this paper we focus on interactive visualization methods based on radial axes (Kandogan, 2000, 2001; Rubio-Sánchez et al., 2017), which map high-dimensional samples onto a two-dimensional space. The transformations are defined through a set of radial axis vectors, each associated with a feature, which users can modify interactively in order to carry out diverse exploratory tasks, such as analyzing correlations, cluster structure, or class separation, or searching for outliers or data with desired characteristics. However, performing feature selection with these methods is cumbersome. On the one hand, a forward selection is impractical, especially for efficiency reasons. On the other hand, while a backwards selection could be implemented with current techniques, the size of the axis vectors and the scale of the plots complicate determining which features should be discarded, from both a visual and an interactive point of view.

Alternatively, in this paper we introduce a new approach based on radial axes that is designed to facilitate performing backwards feature elimination, where users can progressively discard features with a small influence either on the visualizations or on a specific task (e.g., class or cluster separation). Specifically, this is accomplished by employing a set of scaled radial vectors that provide a clearer visual guidance for determining which features have the least impact on the low-dimensional plots, and therefore represent reasonable candidates to be discarded in a backwards elimination process. In practice, analysts determine the contribution of the features to the plots and their related analysis tasks

35 by examining the lengths and orientations of the axis vectors. Moreover, they can also  
36 take into consideration their expertise when deciding whether a feature should belong to  
37 the final selected subset. Lastly, we illustrate the effectiveness of our approach through a  
38 case study related to a real medical chronic conditions data set. Concretely, clinicians have  
39 used the technique, in combination with their expert domain knowledge, in order to ob-  
40 tain insight regarding the discriminative power of the data features for classifying diabetes  
41 and/or high blood pressure patients.

42 The rest of the paper is organized as follows. Section 2 describes the most relevant  
43 methods related to our proposal. In Section 3 we describe our approach based on scaled  
44 axes, illustrating how the proposal can be used to perform visual feature selection. Sec-  
45 tion 4 shows its capabilities through the case study related to medical data. Finally, Sec-  
46 tion 5 presents a discussion with the main benefits and limitations of the proposal, while  
47 Section 6 presents the conclusions and future work.

## 48 2. Related work

49 In this section we present a brief introduction to feature selection methods (with em-  
50 phasis on visual techniques), and describe the most relevant radial axes methods for mul-  
51 tivariate visualization related to our proposal.

### 52 2.1. Feature Selection

53 There is a vast literature on automatic feature selection techniques (Blum & Langley,  
54 1997; Guyon & Elisseeff, 2003; Chandrashekar & Sahin, 2014). *Feature ranking* methods  
55 sort the features according to some criteria and then select the features progressively (*for-*  
56 *ward selection*), consider all of the features initially and discard them sequentially (*back-*  
57 *wards elimination*), or simply apply some threshold to select the top-ranked features. If the  
58 ultimate goal is classification, these strategies are also called *filters*, and discard features  
59 as an independent preprocessing step before training a classifier. Alternatively, *wrapper*  
60 methods select subsets of features according to the accuracy of classification algorithms,  
61 which can be regarded as black boxes that score subsets of features. Lastly, *embedded*  
62 methods use a hybrid strategy that incorporates the feature selection process when training  
63 a particular classifier.

64 The method proposed in this paper can be regarded as a feature ranking procedure  
65 for backwards elimination feature selection. However, instead of defining an automatic  
66 algorithm, it relies on interactive visualizations of data where users can apply their domain  
67 knowledge to steer the process of discarding features. Recently, the data visualization  
68 community has developed several visual feature selection methods and tools that also take  
69 into account user interaction. Most of the approaches propose graphical user interfaces that

70 show several visualizations simultaneously. Some contain well-known graphics in order  
71 to show overviews or properties of the data, while others constitute novel visualization  
72 methods. In order to perform feature selection many of these methods rely on *quality*  
73 *metrics*, which are measures that extract meaningful information about data. While some  
74 of these metrics are popular statistical estimates (correlation, Fisher score, or entropy gain,  
75 among others), many others constitute heuristic measures (May et al., 2011).

76 Several of the earliest proposals are due to Yang et al., which developed hierarchi-  
77 cal methods for visual feature reduction. Yang et al. (2003a) proposes a dimensionality  
78 reduction method based on InterRing visualizations (Yang et al., 2002), which groups  
79 features hierarchically according to their similarity. The method was later extended to  
80 rank and filter out features (Yang et al., 2003b). Guo (2003) describes an interactive tool  
81 using several visualizations (e.g., parallel coordinates (Inselberg & Dimsdale, 1990) and  
82 entropy matrices) to identify subspaces and high-dimensional (hierarchical) clusters. The  
83 approach uses various heuristics, including a measure of the “goodness of a clustering”,  
84 and orderings related to paths on minimal spanning trees (MST). An interactive framework  
85 for ranking features based on ordering histograms and scatter plots is proposed in Seo &  
86 Shneiderman (2005). The work relies on numerous heuristics related to the distributions  
87 that appear in the visualizations (e.g., uniformity, number of outliers or gaps, or modality).  
88 Similarly, Johansson & Johansson (2009) uses heuristics related to the importance of a  
89 feature for correlation, outlier, and cluster detection. By weighting these measures inter-  
90 actively, users can generate feature orderings and reduce the number of features. Ingram  
91 et al. (2010) presents the DimStiller system for feature reduction and analysis. It uses  
92 abstractions (e.g., operators, expressions, or workflows) to combine different visualization  
93 techniques, and structure and guide the data analysis process. In particular, the approach  
94 can be used to determine whether features are meaningful, relationships between features,  
95 or the validity of detected clusters. May et al. (2011) proposes an interactive visualiza-  
96 tion technique denoted as SmartStripes for guiding the feature selection process, which  
97 can be used with categorical features. Tatu et al. (2012) examines clusterings in different  
98 sets of subspaces, which can be interactively explored by relying on subspace similar-  
99 ity and interestingness measures. The visualization tool allows to visualize features and  
100 subsets of features at various levels of detail, through parallel coordinates, lists of scatter  
101 plots, or multidimensional scaling (MDS) (Cox & Cox, 1994) visualizations. Krause et al.  
102 (2014) describes the INFUSE system, which is designed to help interpret how predictive  
103 features are ranked across feature selection algorithms and classifiers. For each feature,  
104 the tool displays a circular glyph depicting information related to several feature selection  
105 methods, which are based on measures of information gain, Fisher score, odds ratios, and  
106 relative risks. In addition, the tool depicts the results of several classification algorithms  
107 for the feature selection methods, across several cross-validation folds. Lastly, Rauber

Method	Task	Reduction approach	Auxiliary visualizations	Quality metric
Yang et al. (2003a)	Dimensionality reduction	Subset selection	InterRing	Similarity
Yang et al. (2003b)	Feature ranking	Subset selection	InterRing	Similarity Importance
Guo (2003)	Feature insight Clustering	Feature reduction Subset selection	Entropy matrix Parallel coordinates Interactive histograms Bar and line charts	Goodness of clustering Maximum conditional entropy MST ordering
Seo & Shneiderman (2005)	Feature ranking	Feature reduction	Score matrix Histograms Scatterplots Box plots	1 and 2-dimensional metrics Modality Outlierness Gaps
Johansson & Johansson (2009)	Feature ranking	Feature reduction	Score matrix Scatter plot matrix Parallel coordinates	Correlation Distribution density
Ingram et al. (2010)	Feature insight Cluster validation	Feature reduction	Scatter plot matrices Correlation matrices Scree plots	Intrinsic dimensionality Variance and correlation MDS stress
May et al. (2011)	Feature insight	Subset selection	Histograms	Mutual information
Tatu et al. (2012)	Clustering	Subset selection	Parallel coordinates Scatterplot lists MDS of subspaces	Subspace redundancy Subspace interestingness
Krause et al. (2014)	Feature insight Classification	Feature reduction Subset selection	Glyphs Bar charts	Information gain Fisher score Odds ratio Relative risk
Rauber et al. (2015)	Classification	Feature reduction	Scatterplots LSP	RFE Random forests

Table 1: Summary of visual feature selection methods in the literature.

108 et al. (2015) proposes a tool for interactive image feature selection including five different  
109 views (observation, projection, feature, group, and feature scoring) that show information  
110 at various levels of detail. The tool uses recursive feature elimination (RFE) (Guyon et al.,  
111 2002) and an ensemble of randomized decision trees (Geurts et al., 2006), and the projec-  
112 tion view employs the least square projection (LSP) (Paulovich et al., 2008) dimensionality  
113 reduction technique.

114 Table 1 presents a brief summary of the previous visual feature selection methods. In  
115 particular, the table considers: (a) the goal or task they are designed for, (b) the reduc-  
116 tion approach, which can consist of progressively discarding features one by one, or of  
117 selecting entire subsets of features in a single step, (c) the auxiliary visualization methods,  
118 and (d), the quality metrics used. It is worth mentioning that the capability of a tool for  
119 feature selection not only depends on the different graphics and the associated interaction  
120 techniques, but also on the nature of the data set, and on the quality metrics used to rank

121 the features (or feature subsets), which are remarkably diverse. Bertini et al. (2011) carries  
 122 out a thorough literature review in order to provide a unified picture of proposed quality  
 123 metrics for high-dimensional data visualization).

## 124 2.2. Radial axes methods

125 In this paper we propose a new approach based on radial axes visualizations that allows  
 126 analysts to perform feature selection effectively. Radial axes methods are popular mul-  
 127 tivariate visualization techniques that produce dimensionality reduction mappings. The  
 128 simplest method is star coordinates (SC) (Kandogan, 2000, 2001), which is an extension  
 129 of the scatterplot for more than two features, and has been used for exploratory tasks such  
 130 as analyzing cluster structure, outliers, or trends. Let  $\mathbf{X}$  be an  $N \times n$  data matrix, con-  
 131 taining  $N$  samples, each characterized by  $n$  features. The method maps high-dimensional  
 132 samples  $\mathbf{x} \in \mathbb{R}^n$  onto a plane by relying on a set of  $n$  axis vectors  $\mathbf{v}_i \in \mathbb{R}^2$ , for  $i = 1, \dots, n$ ,  
 133 with a common origin point. Each  $\mathbf{v}_i$  is associated with the  $i$ -th feature. In particular, the  
 134 low-dimensional representation  $\mathbf{p} \in \mathbb{R}^2$  (also denoted as an “embedded point”) of a sample  
 135  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  is a linear combination of the vectors  $\mathbf{v}_i$ . Formally,

$$\mathbf{p} = x_1 \mathbf{v}_1 + x_2 \mathbf{v}_2 + \dots + x_n \mathbf{v}_n = \mathbf{V}^\top \mathbf{x}, \quad (1)$$

136 where  $\mathbf{V}$  is the  $n \times 2$  matrix whose rows are the vectors  $\mathbf{v}_i$ . The method therefore generates  
 137 linear mappings specified by  $\mathbf{V}$ . In SC, the orientation of an axis vector determines the  
 138 direction in which a feature increases, while the length is related to its contribution to the  
 139 plot. For illustration purposes, Fig. 1(a) shows an example using four features (‘Acceler-  
 140 ation’, ‘Horsepower’, ‘Displacement’, and ‘MPG’) of the Auto MPG data set, available at  
 141 the UCI Machine Learning Repository (Lichman, 2013). The axis vectors have been cho-  
 142 sen to search for cars with large values of ‘Horsepower’ and ‘Acceleration’, but low values  
 143 of ‘MPG’, which would be represented as dots at the top of the plot. The visualization also  
 144 includes an axis vector for ‘Displacement’, which plays a role horizontally. It is important  
 145 to note that although the length of its axis vector is smaller than the remaining lengths,  
 146 its contribution to the plot is important since it has a larger component in the horizontal  
 147 direction.

148 In practice, users can modify the axis vectors interactively in order to carry out diverse  
 149 analysis tasks. However, another possibility is to automatically obtain sets of axis vectors  
 150 from linear methods such as principal component analysis (PCA) (Jolliffe, 2010), inde-  
 151 pendent component analysis (ICA) (Hyvärinen et al., 2001), linear discriminant analysis  
 152 (LDA) (McLachlan, 2004), and so forth. Consider a linear method that maps data points  
 153 onto a plane through  $\mathbf{p} = \mathbf{A}\mathbf{x}$ , where  $\mathbf{A}$  is a known  $2 \times n$  matrix. Clearly, we can build a  
 154 SC model that generates the same plot by setting  $\mathbf{V} = \mathbf{A}^\top$ , due to (1). In other words, we

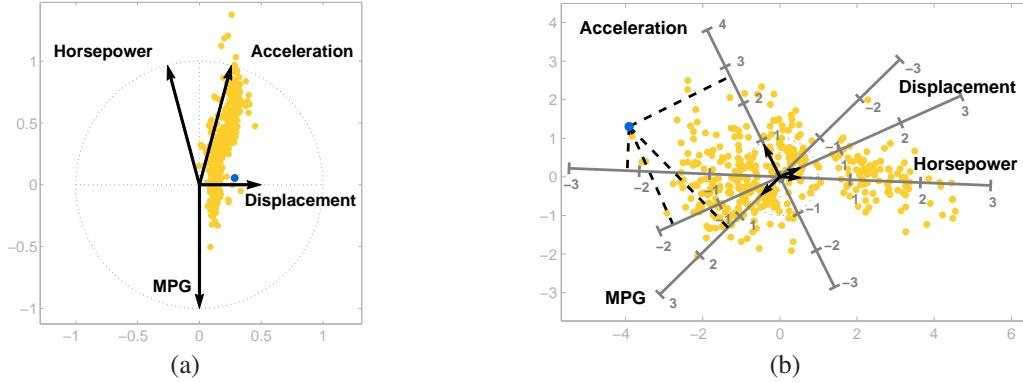


Figure 1: Radial axes plots of the Auto MPG data set: (a) SC plot; (b) ARA plot, where the axis vectors have been selected to generate the PCA projection of the data onto a plane.

155 can recover the SC axis vectors (they would be the columns of  $\mathbf{A}$ ) that lead to the plot related to the linear method. In the SC model, the possibility to visualize these axis vectors, together with the plotted points, allows us to determine relationships between the features and their contribution to the plots. Rubio-Sánchez et al. (2016) introduced this idea to analyze plots based on LDA. Recently, Wang et al. (2017) has denoted it as discriminative star coordinates, and it has also been applied to the results of unsupervised LDA (Ding & Li, 2007), which combines  $k$ -means clustering (MacQueen, 1967) and LDA. Lastly, these works carry out feature selection by only comparing the lengths of the axis vectors. In other words, they do not take advantage of their orientations, which should also be considered (see Section 3.5).

165 Rubio-Sánchez et al. (2017) present a hybrid approach that bridges the gap between SC and principal component biplots (Gabriel, 1971; Gower et al., 2011) called adaptable radial axes (ARA) plots. In SC, users can update the axis vectors freely, but it is difficult to recover high-dimensional data values accurately, which is one of the main disadvantages of the method (Draper et al., 2009). Alternatively, with principal component biplots users can approximate the feature (i.e., data) values of an entire data set as accurately as possible (in a least squares sense) through projections of the embedded points onto ticked axes (see Fig. 1(b)). However, since the axis vectors are fixed in these visualizations, users cannot modify them in order to carry out several exploratory analysis tasks (e.g., searching for data with certain features, or creating different mappings in order to detect outliers or visualize clusters). In ARA plots analysts can update the axis vectors freely, and also approximate data values through projections onto ticked axes. Fig. 1(b) shows an example that uses standardized data. In this case, the means (which are 0) are represented at the origin, and the difference between consecutive tick marks corresponds to one standard deviation of the corresponding feature. Taking this interpretation into consideration, we

180 can approximately determine through orthogonal projections that the car associated with  
 181 the darker blue point (which is also depicted in the SC plot) has a large value of ‘Acceler-  
 182 ation’ (approximately 2.8), and low values of ‘MPG’, ‘Horsepower’ and ‘Displacement’.  
 183 Although the estimated values are simply approximations, it is considerably simpler to ob-  
 184 tain them visually using ticked axes than in the SC graphic (see Rubio-Sánchez & Sanchez  
 185 (2014)). Additionally, it is also worth mentioning that, similarly to SC, it is possible to  
 186 configure the axis vectors to generate any linear mapping. In this example, the particular  
 187 choice of axis vectors leads to a PCA plot of the data.

188 Formally, given a set of axis vectors coded in  $\mathbf{V}$ , ARA plots find the low-dimensional  
 189 embedded point  $\mathbf{p}$  of a data point  $\mathbf{x}$  by solving the following optimization problem:

$$\underset{\mathbf{p} \in \mathbb{R}^2}{\text{minimize}} \quad \|\mathbf{V}\mathbf{p} - \mathbf{x}\|, \quad (2)$$

190 where  $\mathbf{V}\mathbf{p}$  is the vector of approximated values for the data point  $\mathbf{x}$ . Therefore, in ARA  
 191 plots the approximated feature values are the dot products between the embedded points  $\mathbf{p}$   
 192 and the axis vectors  $\mathbf{v}_i$ . In this scenario, the value represented at the endpoint of the axis  
 193 vector is  $\|\mathbf{v}\|^2$ . In addition, a unit of the original feature is located at  $1/\|\mathbf{v}\|$  along the axis,  
 194 which implies that the distance between tick marks separating consecutive integers is also  
 195  $1/\|\mathbf{v}\|$ . Since the length of  $\mathbf{v}$  does not correspond to a unit of a feature (unless  $\|\mathbf{v}\| = 1$ ),  
 196 it cannot be used as a visual reference to indicate the location along the axis where a unit  
 197 would be represented (see Fig. 2(a) for details). Therefore, the method requires drawing  
 198 axis lines together with tick marks representing integers of the features. Without these tick  
 199 marks, users would not be able to approximate data features properly, since it is difficult  
 200 to visually estimate the reciprocal of the length of an axis vector (i.e.,  $1/\|\mathbf{v}\|$ ). Lastly,  
 201 drawing these ticked axes can produce crowded plots even for a small number of features  
 202 (see Section 3.4). The method proposed in this work mitigates this drawback.

### 203 3. Scaled radial axes plots

204 For the purpose of analyzing high-dimensional data and carrying out visual feature  
 205 selection, we propose here a new radial axes method called Scaled Radial Axes (SRA)  
 206 plots. In this section we describe the approach and indicate the main differences with  
 207 other techniques based on radial axes.

#### 208 3.1. Description and mathematical formulation

209 Users in SRA plots will also be able to recover feature values ( $x_i$ ) by relying on or-  
 210 thogonal projections onto axes, similarly to ARA plots. In ARA the approximated values  
 211 correspond to dot products between embedded points and axis vectors, which require axes



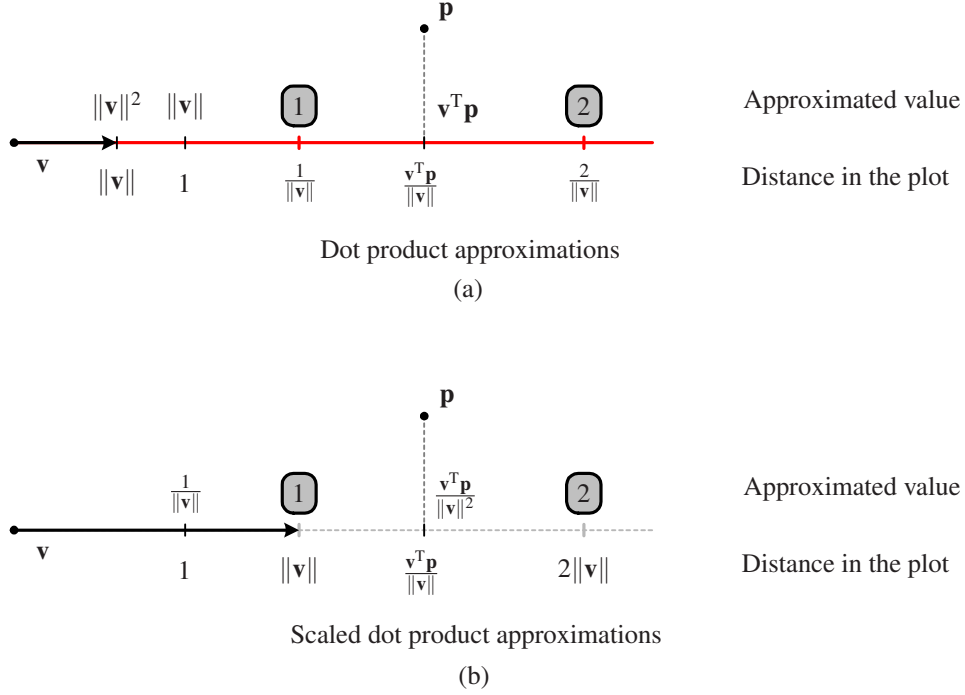


Figure 2: Relationships between approximated values (indicated on the upper part of the horizontal line) and distances in the plots (shown on the lower part of the horizontal line) for: (a) ARA, and (b) SRA. Note that ARA requires axes lines and tick marks (in red) to indicate the values of the approximations.

212 lines and tick marks to indicate the locations associated with integer approximations. Al-  
 213 ternatively, in SRA we consider a more intuitive strategy that uses scaled axes, where a  
 214 unit of a feature is located exactly at the endpoint of its axis vector. Therefore, in this  
 215 scenario the length of an axis vector determines the distance between consecutive integers  
 216 of its corresponding feature. This is illustrated in Fig. 2, which shows the relationships  
 217 between the distances on the plots and the corresponding approximations on the axes, for  
 218 ARA and SRA.

219 In SRA the idea is implemented by recovering the  $i$ -th data feature of a data point  
 220 through the following scaled dot product:

$$\frac{\mathbf{v}_i^T \mathbf{p}}{\|\mathbf{v}_i\|^2}.$$

221 By dividing by the squared Euclidean norm of an axis vector, its endpoint now represents  
 222 a unit of its associated feature, as shown in Fig. 2(b). This allows us to omit drawing  
 223 line axes when the approximations are small (see Section 3.4). Therefore, we define SRA

224 formally through the following optimization problem:

$$\underset{\mathbf{p} \in \mathbb{R}^2}{\text{minimize}} \quad \|\bar{\mathbf{V}}\mathbf{p} - \mathbf{x}\|_2^2, \quad (3)$$

225 where  $\bar{\mathbf{V}}$  is similar to  $\mathbf{V}$ , but in this case each row is divided by its squared norm. Specifi-  
226 cally, the rows of  $\bar{\mathbf{V}}$  are:

$$\bar{\mathbf{v}}_i = \begin{cases} \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2^2} & \text{if } \mathbf{v}_i \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{v}_i = \mathbf{0}. \end{cases} \quad (4)$$

227 The optimal solution to (3) is given by:

$$\mathbf{p} = \bar{\mathbf{V}}^\dagger \mathbf{x}, \quad (5)$$

228 where  $\dagger$  denotes the Moore-Penrose pseudoinverse. The method therefore builds a linear  
229 mapping from the data space onto the observable plane characterized by the matrix  $\bar{\mathbf{V}}^\dagger$ .  
230 We can define the projection of an entire data set in matrix notation through:

$$\mathbf{P} = \mathbf{X}(\bar{\mathbf{V}}^\dagger)^\top, \quad (6)$$

231 where  $\mathbf{P}$  is the  $N \times 2$  matrix whose rows consist of the embedded points. In practice it can  
232 be computed very efficiently, even for large values of  $n$  and  $N$  (see Section 5). Finally,  
233 when  $\bar{\mathbf{V}}$  has full column rank (i.e., when the axis vectors are not all aligned along the same  
234 direction),  $\bar{\mathbf{V}}^\dagger = (\bar{\mathbf{V}}^\top \bar{\mathbf{V}})^{-1} \bar{\mathbf{V}}^\top$ .

### 235 3.2. Influence of the axis vectors on the plots

236 Using  $\bar{\mathbf{V}}$  not only determines how the axes are scaled, but it also affects how the axis  
237 vectors influence the plots, and how users must interact with them. It is important to notice  
238 that shorter vectors will have a stronger impact on the SRA plots, in contrast to longer  
239 vectors when using other radial axes plots. Observe that, when searching for the optimal  
240 embedded point  $\mathbf{p}$ , the optimization problem in (3) naturally focuses on minimizing errors  
241 on shorter axis vectors. In particular, note that the objective function in (3) can be rewritten  
242 as:

$$\sum_{i=1}^n \left( \frac{1}{\|\mathbf{v}_i\|_2^2} \cdot \mathbf{v}_i^\top \mathbf{p} - x_i \right)^2. \quad (7)$$

243 Therefore, if the  $i$ -th axis vector  $\mathbf{v}_i$  is long,  $1/\|\mathbf{v}_i\|_2^2$  will be small and the choice of  $\mathbf{p}$  will  
244 barely affect the  $i$ -th term of the sum in (7). The scaled axis vectors are useful for visual  
245 backwards feature selection since it is easier to spot the longest vectors, associated with  
246 features with a small influence on the plots.

247 However, the length of an axis vector is not the only factor determining the contribution  
 248 of a feature to a plot. To illustrate this, in this work we compute the average displacement  
 249 of the low-dimensional points when a feature is discarded as:

$$f(\mathbf{v}_i) = \frac{1}{N} \sum_{j=1}^N \|\mathbf{p}^{(j)} - \mathbf{q}_{\mathbf{v}_i}^{(j)}\|, \quad (8)$$

250 where  $N$  is the cardinality of the data set,  $\mathbf{p}^{(j)}$  is the embedded point of the  $j$ -th data sample  
 251 for a particular radial axes method, and  $\mathbf{q}_{\mathbf{v}_i}^{(j)}$  is the corresponding low-dimensional point  
 252 when removing the feature associated with the axis vector  $\mathbf{v}_i$ .

253 Fig. 3 shows an example of these average displacements for SC, ARA, and SRA plots.  
 254 Specifically, we generated a random set of  $n = 50$  axis vectors, and a random data set of  
 255  $N = 100$  points. The components of the axis vectors and the values of the data points  
 256 were drawn from a standard normal distribution. Subsequently, we computed the low-  
 257 dimensional points associated with the three methods, and obtained their average displace-  
 258 ments. The dots on the graphics represent pairs  $(\|\mathbf{v}_i\|, f(\mathbf{v}_i))$  and illustrate the average  
 259 displacement of the mapped points when  $\mathbf{v}_i$  is removed from a radial axes plot, as defined  
 260 in (8). The trend for SC and ARA is clearly increasing, but dots do not follow a strictly in-  
 261 creasing pattern as  $\|\mathbf{v}_i\|$  grows. Thus, there are features with longer axis vectors that do not  
 262 contribute as much as others with shorter ones. Similarly,  $f(\mathbf{v}_i)$  does not strictly decrease  
 263 as  $\|\mathbf{v}_i\|$  increases for SRA. For instance, the feature with the second shortest axis vector  
 264 has less impact on the plot than the features with the third to sixth shortest axis vectors.  
 265 Therefore, besides the length of an axis vector, it is necessary to take into account other  
 266 factors such as the orientation of the axis vectors, the arrangement of clusters or classes in  
 267 the plots, or domain knowledge (see Section 3.5). We emphasize this consideration since  
 268 previous works in the literature have only focused on analyzing the lengths of the axis  
 269 vectors.

### 270 3.3. Arbitrary linear mappings

271 Similarly to SC and ARA, it is also possible to select a set of axis vectors in SRA to  
 272 generate any linear mapping from the data space onto the plane. Let  $\mathbf{A}$  be a known  $2 \times n$   
 273 matrix defining the linear transformation to reproduce. Due to (5), we would need to find  
 274 a set of axis vectors for which  $\bar{\mathbf{V}}^\dagger = \mathbf{A}$ . This can be accomplished by first computing the  
 275 pseudoinverse of  $\mathbf{A}$ , which provides  $\bar{\mathbf{V}}$ :

$$\bar{\mathbf{V}} = \mathbf{A}^\dagger, \quad (9)$$

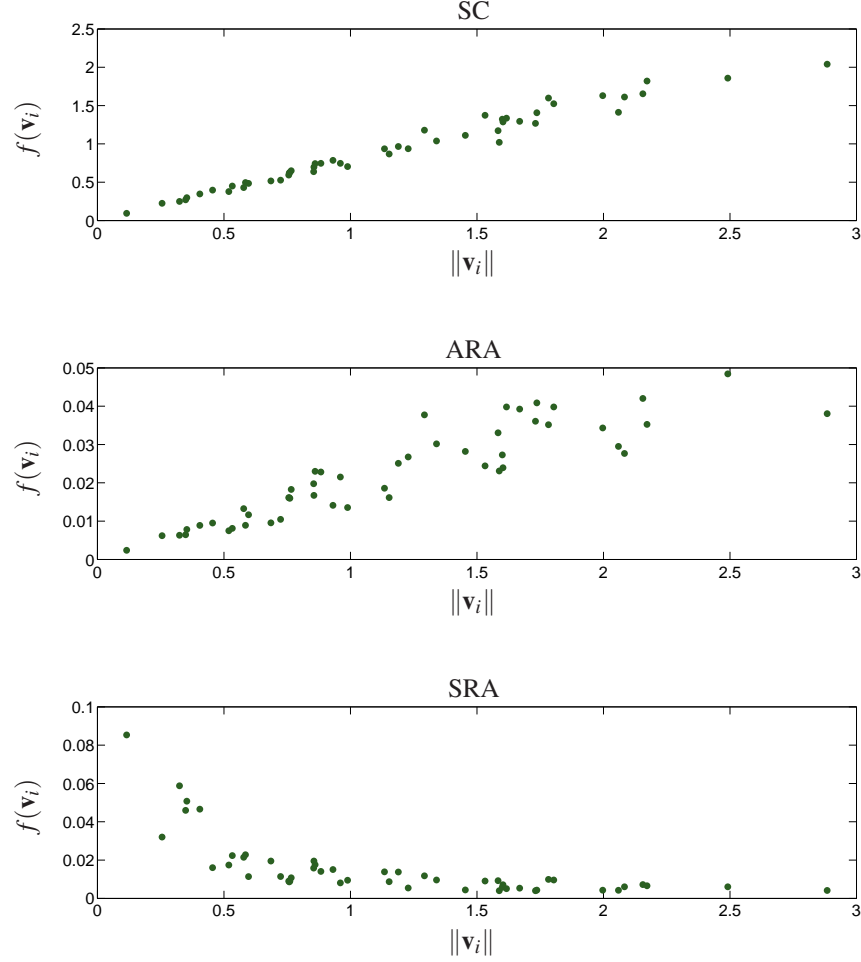


Figure 3: Example of the contribution of axis vectors to plots (in terms of the average displacement of mapped points when removing a feature) depending on their length, for SC, ARA and SRA.

276 since  $\mathbf{M} = (\mathbf{M}^\dagger)^\dagger$  for any matrix  $\mathbf{M}$ . Subsequently, the axis vectors (that form  $\mathbf{V}$ ) can be  
 277 recovered through:

$$\mathbf{v}_i = \begin{cases} \frac{\bar{\mathbf{v}}_i}{\|\bar{\mathbf{v}}_i\|_2} & \text{if } \bar{\mathbf{v}}_i \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \bar{\mathbf{v}}_i = \mathbf{0}, \end{cases} \quad (10)$$

278 which follows from (4), since it defines an involution. The axis vectors are therefore the  
 279 rows of the pseudoinverse of  $\mathbf{A}$ , divided by their squared length. The special case in (10) is  
 280 included by considering that  $\mathbf{A}$  can be any matrix, where some rows of  $\bar{\mathbf{V}}$  could be equal to  
 281  $\mathbf{0}$ . In those cases, the corresponding axes cannot be specified for the features. Thus, their

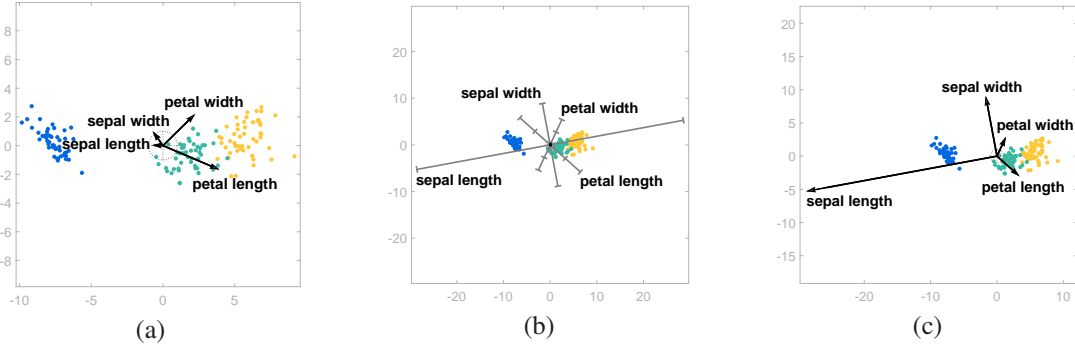


Figure 4: Radial axes plots that produce the LDA mapping of the Iris data set for: (a) SC, (b) ARA, and (c) SRA. The embedded points are colored according to their class. The axis vectors in the ARA plot are very short and are depicted in black near the origin.

282 axis vectors are set to  $\mathbf{0}$ , and the features are ignored when determining the optimal  $\mathbf{p}$ .

283 Fig. 4 shows radial axes plots that produce the LDA mapping of the well-known Iris  
 284 data set (Lichman, 2013). It contains four data features ('petal length', 'petal width',  
 285 'sepal length', 'sepal width') and three classes ('setosa', 'versicolour', 'virginica') that  
 286 identify three species of the iris flower. In particular, we generated the LDA transformation  
 287 automatically (using standardized data) to separate the three classes, and recovered the  
 288 layout of axis vectors that would generate that mapping for SC, ARA, and SRA, in (a), (b),  
 289 and (c), respectively. Note that the plotted points are the same in the three visualizations.  
 290 The SC plot does not incorporate line axes, and therefore users cannot recover feature  
 291 values accurately. The ARA plot mitigates this issue by including ticked axes (but can lead  
 292 to cluttered visualizations for data sets that contain more features). In SRA, the ticked line  
 293 axes are not necessary and the visualization also allows users to recover feature values by  
 294 using the vectors instead of line axes (the endpoints of the vectors indicate the location of  
 295 the units on the axes). Moreover, it is easier to visually identify the less relevant features  
 296 for the class separation task in SRA (longest vectors) than in ARA (shortest vectors), which  
 297 is useful for backwards feature selection. Moreover, in this example the axis vectors in the  
 298 ARA plot are barely visible.

### 299 3.4. Clutter reduction

300 The scaling of the axes is a key contribution regarding the usability of SRA: since  
 301 the vector length visually encodes a unit of the particular feature, it provides the same  
 302 information as the first tick mark on an ARA plot. This allows us to omit drawing line  
 303 axes and their corresponding tick marks when values of the data features are small, which  
 304 reduces clutter considerably.

305 Fig. 5 illustrates an example with the Wine data set available in Lichman (2013). This

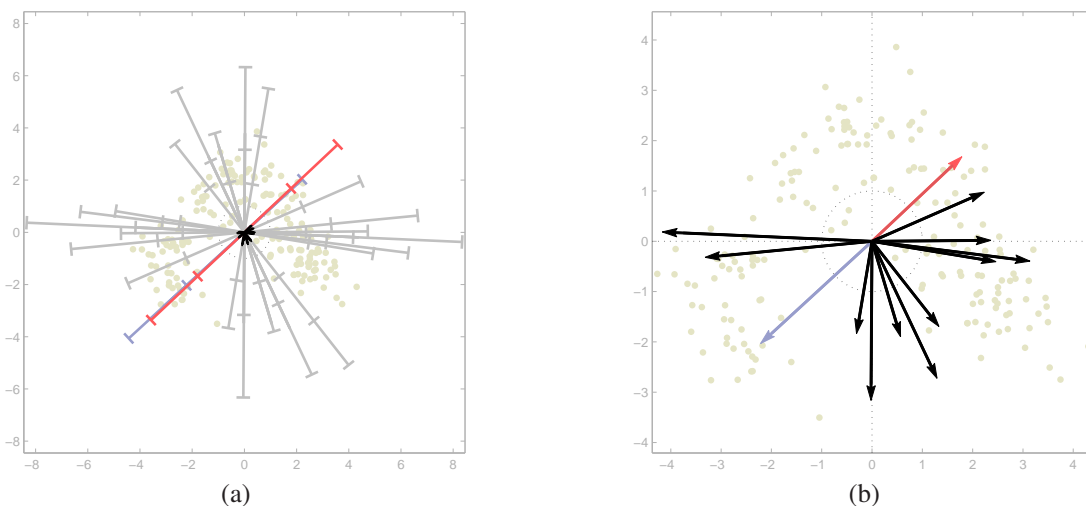


Figure 5: Projection of the Wine data set, composed of 13 features, considering: (a) ARA plot, with axis vectors barely visible due to their small size (depicted in black near the origin), and axes with tick marks; (b) SRA plot using  $\bar{\mathbf{V}}$ , where the axis vectors provide enough visual information to recover original feature values. The clutter reduction when using SRA is apparent (due to the absence of axis lines).

306 data set contains 13 features corresponding to the chemical analysis of three types of wine,  
 307 which we have standardized in a preprocessing stage. The visualization in Figure 5(a) is  
 308 an ARA plot, where we have selected the axis vectors to obtain the PCA projection of  
 309 the data onto a plane. The application of SRA in Fig. 5(b) points out some weaknesses  
 310 of ARA: (1) greater overlap in the ARA plot due to the necessity of drawing the axis  
 311 lines; (2) though the directions of axis vectors are provided by the axis lines, their specific  
 312 orientations are barely visible; and (3) axes can share the same or very similar directions in  
 313 some configurations (e.g., in regular layouts that are often used in the literature), making it  
 314 difficult to distinguish which tick marks are associated with which features. This last issue  
 315 is illustrated in Fig. 5(a), where the colored darker axes exhibit almost identical directions.  
 316 Note that without colors it would not be trivial to identify which tick marks correspond to a  
 317 particular axis. Alternatively, the analogous SRA plot in Fig. 5(b) is less cluttered since it  
 318 does not contain line axes. We have also colored the two vectors that share almost identical  
 319 directions for reference, though this coloring is not necessary in SRA for distinguishing  
 320 the axes and approximating values of the corresponding features. Lastly, when axes are  
 321 omitted it can be easier to incorporate names of features into the plots.

322 In practice, the absence of tick marks in the SRA plot in Fig. 5(b) does not hamper  
 323 users' ability to visually compute projected values severely, in comparison with the radial  
 324 ticked axes plot in Fig. 5(a), which requires them. Note that in radial axes methods the  
 325 features should share a similar scaling, since otherwise features with larger ranges would

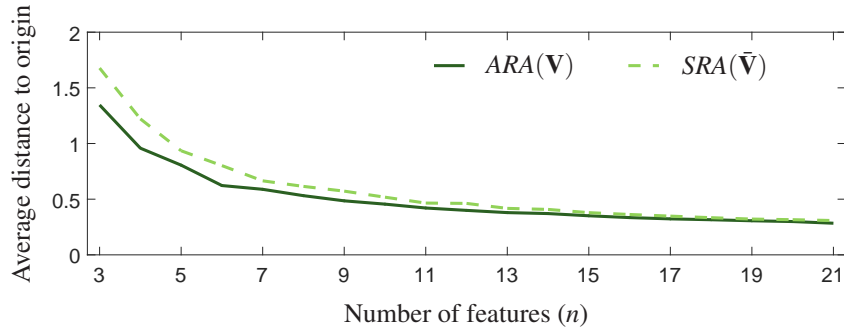


Figure 6: Average distance from embedded points to the origin, for random configurations of vectors and data whose components were drawn from a standard normal distribution.

326 have a greater impact on the resulting plots. Therefore, they are usually standardized,  
 327 transformed to lie in the  $[0,1]$  interval, or centered and normalized to have unit range. In  
 328 all of these cases the absolute values of the approximations corresponding to orthogonal  
 329 projections onto the axes are generally lower than two. Therefore, users can approximate  
 330 these values accurately by relying exclusively on the depicted axis vectors, whose end-  
 331 points are equivalent to one tick mark in a ticked axis.

332 Furthermore, the projections onto the axes in SRA are small not only because the data  
 333 are standardized, but also due to the clumping effect of the projections, which tends to  
 334 map points closer to the origin as the number of features increases. This effect is shown in  
 335 Fig. 6, which shows average distances from embedded points to the origin as a function of  
 336 the number of features ( $n$ ). The results were averaged over 200 trials of random configu-  
 337 rations of vectors, where we mapped 50 samples in each trial. The components of the axis  
 338 vectors, and the values of the data points, were drawn from a standard normal distribution.

339 Finally, standardization has two main benefits. Firstly, a unit of a feature represents  
 340 one standard deviation. Thus, the length of an axis vector in SRA, or the location of the  
 341 first tick mark in ARA, have a clear statistical meaning. This is important to simplify the  
 342 graphics, since it allows us to omit numerical labels next to the tick marks (see Fig. 1(b)).  
 343 Secondly, Rubio-Sánchez & Sanchez (2014) showed that the approximations are more  
 344 accurate when the data are centered.

### 345 3.5. Interactive visual feature selection for class separation

346 Since the scaling introduced in SRA highlights the least important features, the tech-  
 347 nique is appropriate for visual sequential backwards feature selection. In practice, users  
 348 can eliminate features progressively by considering their contribution to a specific plot,  
 349 which is affected by the lengths and directions of the axis vectors. They can also decide to  
 350 maintain or discard features according to their domain knowledge.

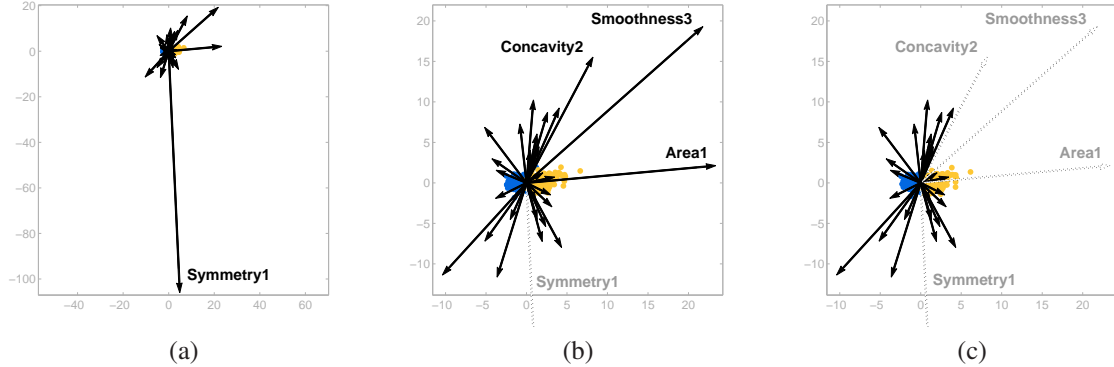


Figure 7: Interactive visual feature selection. SRA plots related to LMNN for the Breast Cancer Wisconsin Diagnostic data set: (a) considering all features, (b) after removing the ‘Symmetry1’ feature; and (c) when removing features named ‘Smoothness3’, ‘Area1’, and ‘Concavity2’.

351 In addition, assuming the data are categorized into several classes, it is possible to  
 352 recover the axis vectors in SRA to generate plots related to linear methods designed to  
 353 enhance classification performance. The most popular linear method is LDA, which max-  
 354 imizes the ratio between the inter-class and intra-class variance. In this paper we will also  
 355 rely on metric learning approaches such as large margin nearest neighbor (LMNN) (Wein-  
 356 berger & Saul, 2009), and neighbourhood components analysis (NCA) (Goldberger et al.,  
 357 2005), whose goal consists of enhancing nearest neighbor classification. The resulting  
 358 SRA plots will provide insight regarding the less discriminative features in the data.

359 For instance, Fig. 7 shows an SRA plot associated with a LMNN mapping of the  
 360 Breast Cancer Wisconsin Diagnostic data set (Alcala-Fdez et al., 2008), which includes  
 361 30 features from a digitized image of a fine needle aspirate of breast mass, used to de-  
 362 termine if a tumor is benign (darker blue dots) or malignant (lighter orange dots). The  
 363 data set contains information regarding 10 characteristics (radius, texture, perimeter, area,  
 364 smoothness, compactness, concavity, concave points, symmetry, and fractal dimension)  
 365 of the cell nuclei present in the image. For each characteristic the data set includes three  
 366 types of measurements: (1) mean, (2) standard error, and (3) the mean just considering the  
 367 three largest values for each image. In the plots we have appended a numerical suffix to  
 368 the names of the features to indicate the type of measurement. Fig. 7(a) shows an SRA plot  
 369 when using the 30 features of the data set. In contrast to SC or ARA plots, features with  
 370 long vectors can be easily detected in SRA, and discarded in a backwards feature selection  
 371 process. In this case, the axis vector for ‘Symmetry1’ is clearly larger than the rest. This  
 372 implies that it barely affects the plot, and it is likely the least discriminative feature. After  
 373 discarding ‘Symmetry1’, the SRA plot is shown in Fig. 7(b), where axis vectors related to



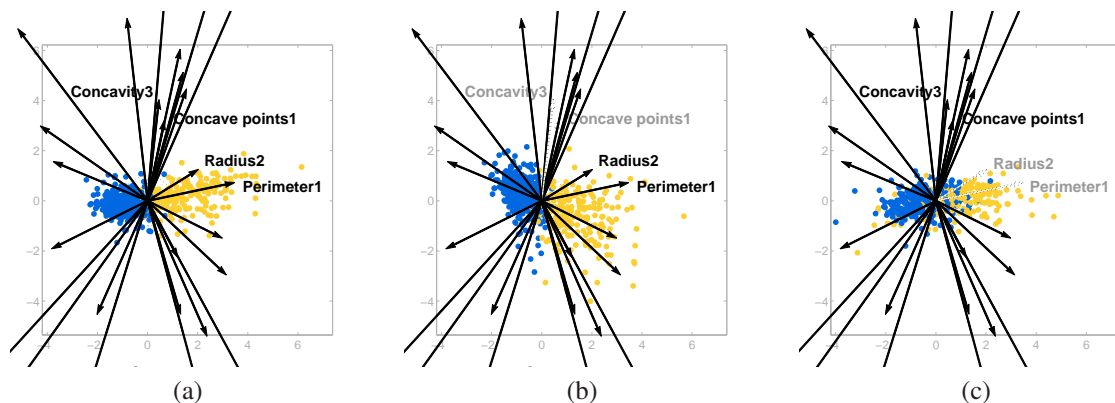


Figure 8: SRA plots related to LMNN for the Breast Cancer Wisconsin Diagnostic data set: (a) zoom of Fig. 7(c); (b) effect of removing the ‘Concave points1’ and ‘Concavity3’ features in (a); (c) effect of discarding ‘Radius2’ and ‘Perimeter1’ in (a).

374 ‘Smoothness3’, ‘Area1’, and ‘Concavity2’ are also longer than the rest. Thus, we can also  
 375 omit these features by focusing on the lengths of the axis vectors, assuming it is appropri-  
 376 ate according to domain knowledge. The resulting plot is shown in Fig. 7(c), where the  
 377 locations of the points are very similar to those in Fig. 7(b).

378 As previously indicated, the direction of an axis vector also constitutes a key factor  
 379 regarding the importance of a feature in a plot. Note that the low-dimensional points will  
 380 move roughly in the direction of an axis vector when the corresponding feature is removed.  
 381 Thus, for separating classes (or clusters) in the two-dimensional plot, we can also discard  
 382 features whose axis vectors are roughly perpendicular to the direction separating these  
 383 classes, even if those axis vectors are short. Fig. 8 illustrates this idea. In particular,  
 384 Fig. 8(a) is just a zoomed version of the plot in Fig. 7(c), where both classes are separated  
 385 fairly well horizontally. Observe that there are several axis vectors whose orientations  
 386 are roughly perpendicular to the class separation direction. Therefore, although omitting  
 387 them could originate large displacements of the plotted points, the two classes should  
 388 remain fairly separated. Specifically, in the plot in Fig. 8(b) we have removed the features  
 389 ‘Concave points1’ and ‘Concavity3’, which have relatively short axis vectors. The low-  
 390 dimensional points therefore move vertically, but this barely alters the overlap between  
 391 classes. Instead, in Fig. 8(c) we have eliminated ‘Radius2’ and ‘Perimeter1’, since their  
 392 axis vectors point in the separation direction. In this case, although their lengths are similar  
 393 to those for ‘Concave points1’ and ‘Concavity3’, the points move roughly horizontally.  
 394 This substantially increases the overlap between the classes, which indicates that these  
 395 features should belong to the final feature subset.

396 The process can continue by considering the lengths and orientations of other axis

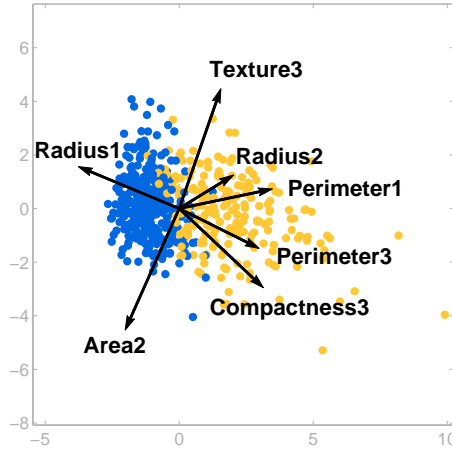


Figure 9: SRA plot illustrating class separation after selecting seven out of the thirty features of the Breast Cancer Wisconsin Diagnostic data set.

397 vectors (and possible domain knowledge), and by analyzing the class separation in the two-  
 398 dimensional plots. The idea is to obtain a final subset of features that allows to separate  
 399 classes reasonably well. Fig. 9 shows an example of an SRA plot where we have retained  
 400 seven of the original thirty features of the Breast Cancer Wisconsin Diagnostic data set.

401 Lastly, we measure the quality of SRA projections for class separation as carried out in  
 402 Leban et al. (2006), by computing the leave-one-out accuracy of a voting  $k$ -nearest neigh-  
 403 bor ( $k$ -nn) classifier (Duda et al., 2001) applied on the plotted two-dimensional points.  
 404 Specifically, we used  $k = \sqrt{N}$ , where  $N$  is the number of samples in the data set, as sug-  
 405 gested by Dasarathy (1991). Thus, for the Breast Cancer Wisconsin Diagnostic data set we  
 406 chose  $k = 24$ , since it contains  $N = 569$  samples. We obtained a quality of class separation  
 407 of 96.66% when considering the plot in Fig. 7(a) that involves all of the 30 features in the  
 408 data set. The score only dropped to 93.32% when considering the plot in Fig. 9, which  
 409 uses the reduced set of seven features.

#### 410 4. Case study: analyzing chronic conditions

411 In this section we describe a case study in which clinicians used SRA for visual feature  
 412 selection related to chronic conditions.

##### 413 4.1. Chronic conditions fundamentals

414 Chronic diseases constitute a well-known problem in current societies, mainly due to  
 415 the major demographic changes throughout the world over the past few years. On the one

416 hand, the percentage of people over 65 years of age is expected to increase in developed  
417 regions (McNicoll, 2002). On the other hand, it is estimated that by the year 2050 about  
418 20% of the whole world population will exceed 65 years. There are also clear positive  
419 correlations between age, chronic conditions, and the use of health services. According  
420 to Organization et al. (2005), chronic diseases account for 60% of global deaths, and trig-  
421 ger 75% of public health expenditure. Therefore, it is important to determine the diseases  
422 that present the highest prevalence, and to identify the factors that best characterize them.

423 Two diseases that highly contribute to the complex chronic patient group are diabetes  
424 mellitus (DM) and high blood pressure (HBP, also called essential arterial hypertension).  
425 Not only are they notoriously widespread, but their frequency increases with age, and pa-  
426 tients maintain their chronic condition until their death. Specifically, DM is one of the  
427 leading chronic diseases in developed countries. It entails many consequences, both from  
428 a clinical and social viewpoint, since it increases the risk of many serious health prob-  
429 lems. For example, vascular disease is the diabetes complication that can have a more  
430 severe prognosis, since it can be accompanied by damage to the coronary arteries, which  
431 may lead to myocardial infarction or limb amputation. Other complications of diabetes  
432 include kidney problems and blindness. HBP, which is diagnosed when diastolic/systolic  
433 blood pressure is 140/90 mmHg or greater, appears among 18% of those who suffer from  
434 chronic conditions (Organization, 1999). It can be associated with the onset of other med-  
435 ical conditions such as chronic kidney disease, and it is also related to DM. The simulta-  
436 neous presence of chronic diseases (comorbidities) can have dramatic consequences. For  
437 instance, HPB in patients with DM raises the risk of cardiovascular disease.

#### 438 4.2. *Chronic conditions data*

439 In this case study we used data provided by Hospital Universitario de Fuenlabrada  
440 (HUF) in Madrid, Spain. In order to identify patients with certain chronic diseases, a  
441 Patient Classification System (PCS) was applied. In essence, a PCS is a medical decision  
442 tree with clinically validated rules, which groups patients according to their health status  
443 and resource consumption. Berlinguet et al. (2005) analyzed different PCS and concluded  
444 that the so-called Clinical Risk Groups (CRGs) offered the best performance according to  
445 three criteria: clinical relevance of the grouping, resource prediction, and ease of use. This  
446 was the reason for using the CRGs (Averill et al., 1999; Hughes et al., 2004) to determine  
447 a patient's health status. CRGs are hierarchically organized into nine core categories, from  
448 CRG-1 (healthy user) to CRG-9 (catastrophic).

449 Our data set contains information relative to demographic features (age and gender),  
450 diagnoses from primary and specialized care centers, and pharmaceutical drug dispen-  
451 sation during one year. Diagnoses were coded by considering three digits, as stated in  
452 the International Classification of Diseases, 9th revision, Clinical Modification (ICD-9-

453 CM) (Centers for Disease Control and Prevention, 2011). Medical drugs were specified  
454 through five characters, according to the Anatomical Therapeutic Chemical (ATC) Classi-  
455 fication System (Norwegian Institute of Public Health, 2017) used in Europe. CRGs used  
456 this information to assign each patient to a single mutually exclusive health status or risk  
457 group.

458 In this paper we analyzed three chronic conditions (i.e., categories): crg-5192 (HBP),  
459 crg-5424 (DM), and crg-6144 (DM and HBP). The first digit of the CRG-code refers to the  
460 core group, while the next three digits are associated with the chronic condition category.  
461 Specifically, HUF provided us with data of 17792 patients associated with the three chronic  
462 statuses of interest during the year 2012: 12447 for crg-5192, 2166 for crg-5424, and 3179  
463 for crg-6144. Since class-imbalance is a well-known issue in medical research (Soguero-  
464 Ruiz et al., 2016; Fernández-Sánchez et al., 2017), we adopted an undersampling strategy  
465 taking into account the size of the minority group. Thus, we randomly selected 2166  
466 patients from each group.

467 In a previous study we performed a descriptive analysis of diagnosis codes and demo-  
468 graphic features in the group of only chronic hypertensive patients (Fernández-Sánchez  
469 et al., 2017). Regarding the features in the current work, we have also considered medical  
470 drugs apart from diagnosis. Each code of diagnosis and medical drug has been considered  
471 as a different feature. In particular, each patient is described by a total of 1517 features  
472 for diagnoses, and 746 for medical drugs. The features are integers that count the number  
473 of times that a particular patient has been diagnosed with a certain condition, or has been  
474 dispensed a particular drug. Around half of the features had a zero count for every single  
475 patient, and were therefore discarded. In addition, we reduced the data set even further by  
476 computing the entropy gain of each feature according to Rauber & Steiger-Garção (1993),  
477 and by selecting the 50 features with the highest gain. According to the domain knowl-  
478 edge of the clinicians who participated in the case study, the resulting subset of features  
479 contained the most relevant features related to the chronic conditions under study.

#### 480 4.3. *Visual feature selection with SRA*

481 Since the dimensionality of the data (50) is still high, further feature selection proce-  
482 dures can be useful for identifying features with a greater clinical relevance for character-  
483 izing the chronic conditions. In our case study, the medical doctors used SRA, coupled  
484 with linear methods for classification, as a basis for performing a sequential backwards  
485 visual feature selection. Specifically, the goal was to determine which features were more  
486 helpful for discriminating between health statuses: (i) HBP, (ii) DM, and (iii) HBP and  
487 DM. Therefore, the clinicians used SRA to graphically identify different health groups,  
488 and to evaluate or confirm (in consonance with domain knowledge) the impact of each  
489 feature on the plots designed for class separation. Since clinicians were not experts in data

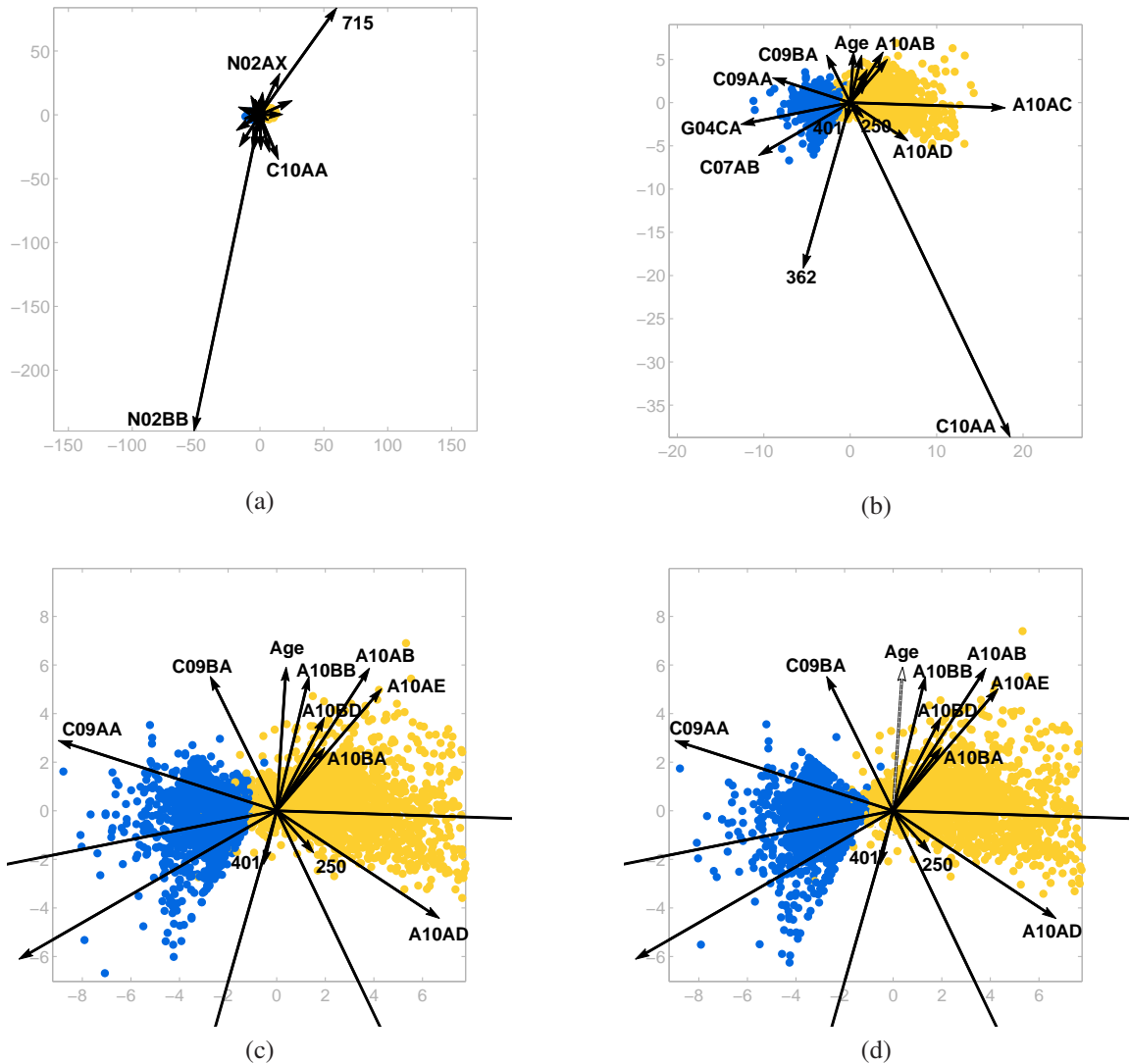


Figure 10: SRA plots related to LMNN for patients with hypertension (darker points, crg-5192) and diabetes (lighter points, crg-5424) considering 50 and 16 features, in (a) and (b), respectively. The plot in (c) represents a zoom of (b), and in (d) we show the (minor) effect of removing the feature ‘Age’.

490 visualization methods, we provided explanations of the main properties of SRA, as well  
 491 as assistance throughout the process.

492 Firstly, the medical doctors analyzed which features contributed more to distinguishing  
 493 between the hypertensive and diabetic groups (crg-5192 vs. crg-5424). This is the simplest  
 494 scenario when considering chronic conditions, since the health statuses are characterized

495 by only one chronic condition. Fig. 10 shows SRA plots associated with the LMNN map-  
496 ping of the (standardized) data set, where the lighter (yellow) and darker (blue) points  
497 represent patients with DM, and HBP, respectively. In Fig. 10(a) we used the initial 50  
498 features. The clinicians then progressively discarded features by relying on the visualiza-  
499 tions and their own expertise until obtaining the plot in Fig. 10(b), which only contains  
500 16 features. The quality of class separation only decreased from a score of 98.66% (when  
501 using the initial 50 features) to 98.61% when considering just 16 features (in this case we  
502 used the voting 66 – *nn* classifier, since there are  $N = 4332$  samples).

503 The plot in Fig. 10(c) is simply a zoom of Fig. 10(b), where we can gain insight regard-  
504 ing the most relevant features for classifying patients with a single chronic condition. In  
505 this example, these features are mainly those oriented horizontally, since classes are sep-  
506 arated along that direction. For instance, the features related to the drug codes ‘G04CA’  
507 (alpha-adrenoreceptor antagonists) and ‘C09AA’ (angiotensin-converting-enzyme inhibitors,  
508 plain) point towards the *crg-5192* class, as expected by the clinicians. Analogously, sev-  
509 eral axis vectors are oriented towards the *crg-5424* class. Their contribution to the plots, as  
510 suggested by their lengths and orientations, was in accordance with the clinician’s back-  
511 ground knowledge. For example, the axis vectors for drug codes ‘A10AB’ (insulins and  
512 analogues for injection, fast-acting), ‘A10AE’ (insulins and analogues for injection, long-  
513 acting), ‘A10BA’ (biguanides), or ‘A10BD’ (combinations of oral blood glucose lowering  
514 drugs) all have positive components along the plot’s X axis, since they point towards the  
515 first quadrant. Thus, they are clearly related to diabetes. The feature for the diagnosis  
516 code ‘250’ (DM) also appears pointing towards the diabetic group, and has a higher con-  
517 tribution than the ATC codes, since its axis vector is shorter. Clinicians also suggested to  
518 retain the drug code ‘C10AA’ (HMG CoA reductase inhibitors) in spite of the long length  
519 of its axis vector, since it could have some relation with diabetic patients. Finally, regard-  
520 ing the ‘Age’ feature, the length of its axis vector is similar to that of the remaining ones.  
521 However, it does not play a key role in separating the *crg-5192* and *crg-5424* groups, since  
522 its axis vector is roughly perpendicular to the direction that separates the classes. This  
523 also occurs for other features like the diagnosis code ‘401’ (essential hypertension). If the  
524 ‘Age’ feature is removed (as shown in Fig. 10(d)), the classes remain clearly separated,  
525 and the quality of class separation is enhanced to 99.01%.

526 For comparison purposes, in Fig. 11 we show SC and ARA plots related to the LMNN  
527 mapping, with the initial 50 features. In both cases shorter vectors have a weaker impact  
528 on the resulting plots. Thus, in practice it is required to zoom in several times to be able  
529 to identify the features to be removed. In the example, the initial SC plot is shown in  
530 (a), while (b) and (c) show 4x and 40x zooms, respectively. Similarly, (d) is the initial  
531 ARA plot, while (e) and (f) show 20x and 100x zooms, respectively. Observe that the  
532 axis vectors (and the axis lines in ARA) overlap considerably, which makes it difficult

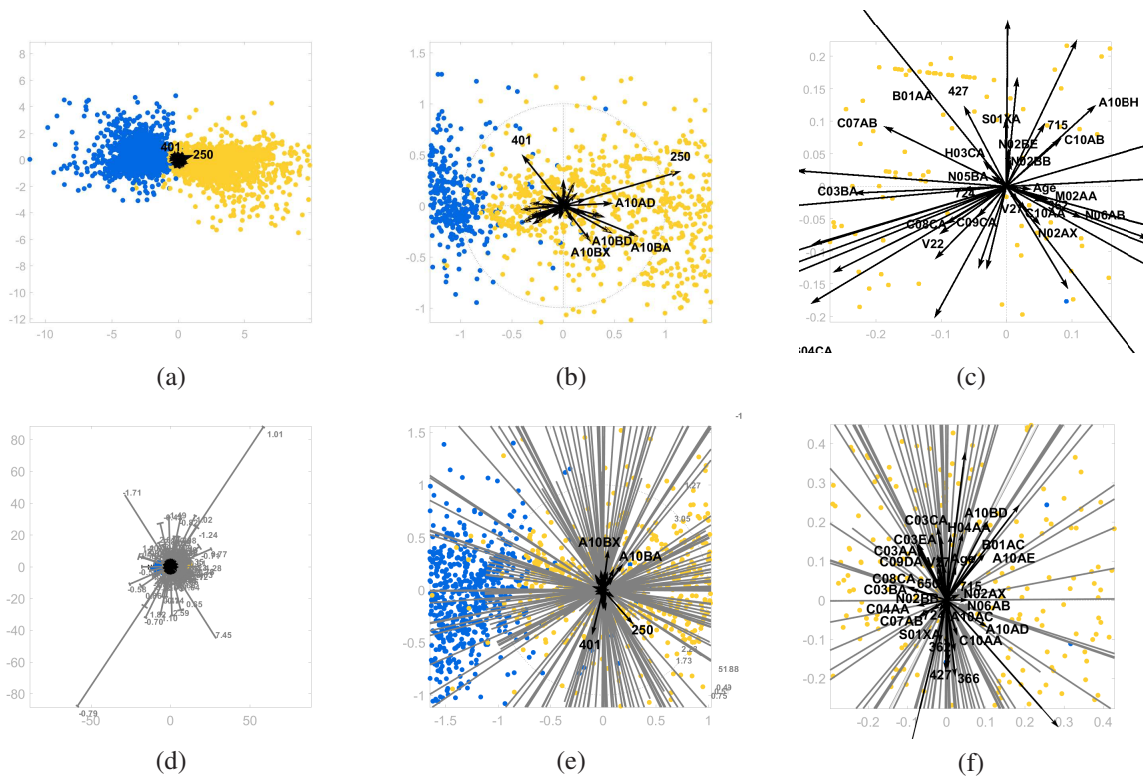


Figure 11: SC and ARA plots related to the SRA plot in Fig. 10 (a) with 50 features. The initial configuration of the SC plot is shown in (a), while (b) and (c) show 4x and 40x zooms, respectively. Analogously, (d) contains the initial ARA plot, while (e) and (f) show 20x and 100x zooms, respectively. On the one hand the axis vectors (and axes lines) overlap considerably. On the other hand, we can lose the distribution of the plotted points when zooming.

533 to visualize and select the shortest axis vectors. In addition, depending on the scale of  
 534 the data, the projected points may fall outside of the plot. Thus, we can lose the overall  
 535 picture of the data set, which is necessary for considering the orientations of the vectors  
 536 (in this case, the direction that separates the classes). In our experiments, all clinicians  
 537 were able to immediately obtain the longest axis (‘N02BB’) using SRA, and agreed to  
 538 remove it (see Fig. 10(a)). However, when using SC and ARA they had to zoom in several  
 539 times, obtaining the plots in (c) and (f), before deciding on the least relevant features.  
 540 Most importantly, they did not agree on the feature to be removed, as some vectors were  
 541 of similar size.

542 In the next study the data set was expanded by including a third health status encom-  
 543 passing both chronic conditions, diabetes and hypertension (crg-6144). In this case, we  
 544 selected a total of 6498 patients (2166 of each health status), and tested our approach by

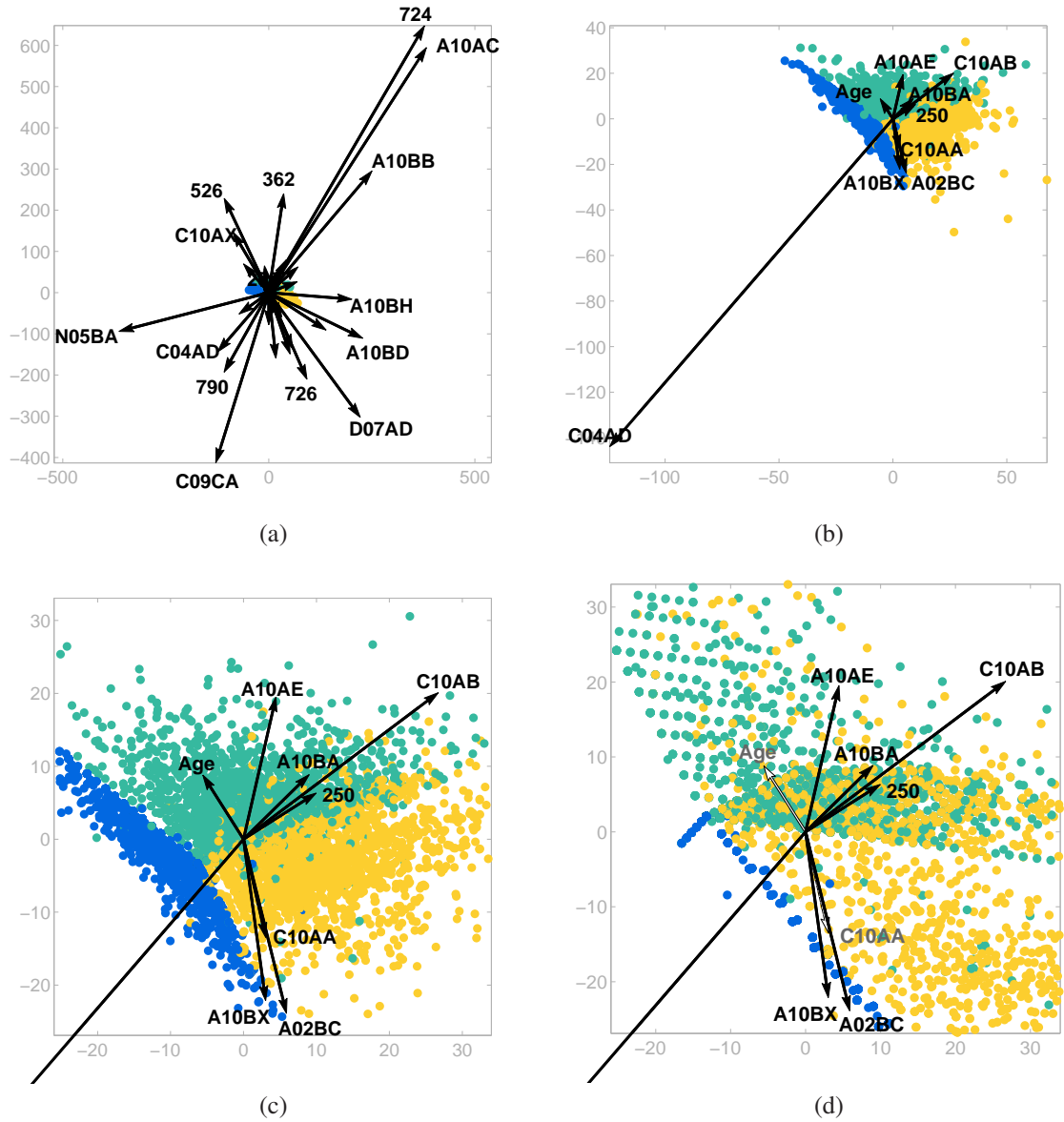


Figure 12: SRA plots related to NCA for patients with just hypertension (darker blue points, crg-5192), just diabetes (lighter orange points, crg-5424), and both comorbidities (mid-range green color, crg-6144) considering 50 and 9 features, in (a) and (b), respectively. The plot in (c) is a zoom of (b), and in (d) we show the (strong) effect of removing the features ‘Age’ and ‘C10AA’.



545 relying on the NCA mapping of the data set. Fig. 12 shows several SRA plots associ-  
546 ated with NCA, where the lighter (yellow), darker (blue), and mid-color (green) points  
547 represent patients with DM (crg-5424), HBP (crg-5192) and both chronic conditions (crg-  
548 6144), respectively. Similarly to the first study, we generated an initial plot by using all  
549 of the 50 features, as shown in Fig. 12(a). The quality of class separation according to a  
550 nearest neighbor classifier was 92.67% (we used  $k = 81$ , since  $N = 6498$ ). Subsequently,  
551 the clinicians progressively eliminated features by relying on the visualization and their  
552 domain knowledge until obtaining the plot in Fig. 12(b), which only contains 9 features  
553 and provides a quality of class separation of 87.17%.

554 We can observe the axis vectors (and their contribution) more clearly in Fig. 12(c),  
555 which is a zoom of Fig. 12(b). On this occasion, clinicians did not select the diagno-  
556 sis code '401' because there were other features with more influence for separating both  
557 groups. Instead, although in the first study the drug code 'C10AA' (HMG CoA reductase  
558 inhibitors) did not contribute much in distinguishing between hypertensive and diabetic  
559 patients (according to the layout of vectors obtained when reproducing LMNN), the clin-  
560 icians suggested to retain it since in their opinion it had a clear relation to diabetes. In  
561 this case, it is apparent that the feature 'C10AA' is key for separating the groups (note  
562 that its axis vector is one of the shortest ones). This confirms the medical knowledge that  
563 reductase inhibitors are related to diabetic patients. Likewise, the feature 'Age' does have  
564 a strong impact on class separation, since individuals in CRGs with chronic comorbidi-  
565 ties (crg-6144) tend to be older than patients with just one chronic condition (crg-5192 or  
566 crg-5424). 'Age' is especially relevant for patients with diabetes, which supports existing  
567 knowledge about juvenile diabetes. Finally, in order to visually confirm the importance of  
568 both features ('C10AA' and 'Age') we discarded their axis vectors. The resulting plot is  
569 shown in Fig. 12(d), where the lighter (crg-5244) and mid-color (crg-6144) classes clearly  
570 overlap. In this case, the quality of class separation dropped to 75.45%.

571 The study carried out, involving clinicians and a real medical data set, shows that  
572 SRA can be a valid tool when it is used by domain experts without previous experience  
573 in interactive visual data analysis tools. The visualizations have allowed the clinicians at  
574 HUF to confirm previous medical knowledge, and to obtain new insight into the area of  
575 application.

## 576 5. Discussion

577 In practice, analysts can use radial axes plots for visual feature selection by studying  
578 the impact of the features on a plot. However, it is problematic to use these visualizations  
579 in a sequential forward selection process, mainly due to the large number of plots that  
580 users would have to analyze. Note that having a subset of  $m < n$  features, it would be

581 necessary to visualize the  $n - m$  additional plots that include one more feature in order to  
582 expand the subset. Since this procedure would be carried out multiple times, the number of  
583 visualizations would be excessive in a practical setting. In particular, this approach would  
584 require  $(m + 1)(n - m/2)$  visualizations for obtaining a subset of  $m$  features. Alternatively,  
585 users in a sequential backwards elimination procedure analyze a single plot to discard one  
586 of the features. Thus, this approach requires analyzing  $n - m$  visualizations in order to  
587 choose a subset of  $m$  features, which is much smaller than the number required by the  
588 sequential forward selection scheme. Thus, if  $m$  is some percentage of  $n$  (i.e.,  $\alpha n = m$ , with  
589  $\alpha \in (0, 1)$ ), then the forward selection strategy requires on the order of  $n^2$  visualizations,  
590 while the backwards approach needs on the order of  $n$  plots. Moreover, when performing  
591 a backwards selection it is also possible to identify an entire group (i.e., set) of features to  
592 discard by analyzing a single plot, which can speed up the selection process notably when  
593 the initial number of features is large.

594 In SRA a backwards feature selection is implemented by removing longer axis vec-  
595 tors, which are easy to spot. In SC and ARA it is possible to perform a similar feature  
596 elimination by discarding shorter axis vectors. However, as shown in Fig. 11, it is more  
597 difficult to identify these axis vectors. In practice, analysts may need to zoom in on the  
598 plots considerably, which is not only time-consuming, but the overall view of the data can  
599 be lost in the resulting graphic, since many of the projected points may not appear in the  
600 plot. Therefore, in SC and ARA it can be harder to take advantage of the directions of the  
601 axis vectors.

602 Although methods based on radial axes can represent as many variables as desired, in  
603 practice  $n$  is usually small (see (Gabriel, 1971; Kandogan, 2000, 2001; Chen & Liu, 2004;  
604 Zhang et al., 2006; Tsai & Chiu, 2008; Sun et al., 2008)). Note that if  $n$  is large a feature  
605 reduction process would be time-consuming and cumbersome, mainly due to the overlap  
606 between the axis vectors. In that case one solution consists of carrying out a preliminary  
607 feature reduction with an automatic method (in Section 4.2 we have used the entropy gain  
608 to reduce the number of features). Another possibility is to generate an SRA plot and  
609 eliminate the features related to long axis vectors, according to a length threshold, or to a  
610 particular number of features the analysts may wish to retain before applying the proposed  
611 feature reduction approach. Another limitation of the approach is related to the type of  
612 data it can support. In particular, all of the radial axes methods described in this paper  
613 require using numerical data (it is possible to use binary features).

614 In order to evaluate the method's potential for data analysis, we have developed a  
615 data visualization prototype in MATLAB<sup>®</sup> using the toolbox for dimensionality reduc-  
616 tion (Maaten, 2015). In preliminary usability tests, users were able to carry out: i) tasks  
617 directly related to the technique like classification, clustering, feature selection, outlier  
618 detection, or attribute value estimation; and ii) other basic data analysis tasks like those

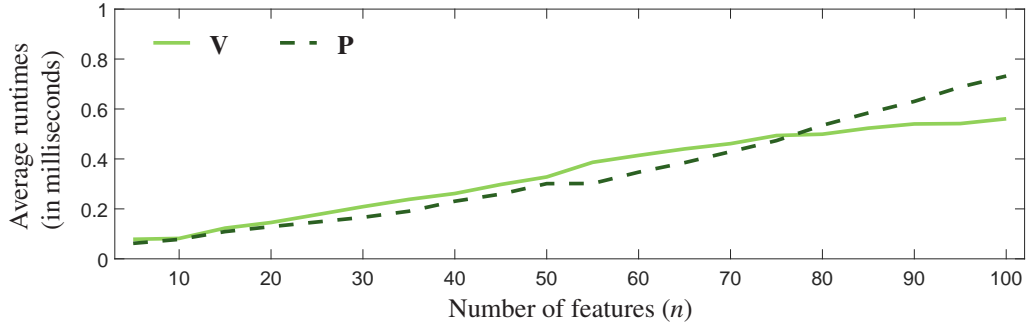


Figure 13: Average runtimes for computing the axis vectors ( $\mathbf{V}$ ) given some initial linear transformation matrix through (9) and (10), and for calculating 10000 embedded points ( $\mathbf{P}$ ) through (6).

619 described in Amar et al. (2005) and Yi et al. (2007), such as retrieving values, determining  
 620 correlations, filtering, etc.

621 Regarding the efficiency of the approach, it is worth mentioning that the key factor  
 622 depends on the computational cost of the chosen linear method (e.g., LDA, LMNN, NCA,  
 623 etc.), which provides a particular  $2 \times n$  matrix  $\mathbf{A}$ . The process of determining the axis  
 624 vectors  $\mathbf{V}$  through (9) and (10), as well as computing the embedded points ( $\mathbf{P}$ ) through (6)  
 625 can be carried out in the order of microseconds, even for a large number of features ( $n$ ),  
 626 since these operations can be carried out in linear time with respect to  $n$ . Figure 13 shows  
 627 average runtimes needed to compute  $\mathbf{V}$  given some random initial matrix  $\mathbf{A}$ , and to project  
 628  $N = 10000$  random high-dimensional points ( $\mathbf{X}$ ), for several values of  $n$ . The results were  
 629 averaged over 1000 trials, and the components of  $\mathbf{A}$  and  $\mathbf{X}$  were drawn from a standard  
 630 normal distribution. In particular, the simulation was carried out on a personal computer  
 631 with a fourth generation Intel<sup>®</sup> Core<sup>™</sup> i7-4712HQ 3.3 GHz processor and 16 GB of RAM.  
 632 It is apparent that the calculations can be carried out in real time.

633 Finally, the proposed visualization method is an exploratory data analysis tool that  
 634 can lead to interesting and possibly unexpected discoveries in an overview phase of a  
 635 data mining process (Shneiderman, 1996; Witten & Frank, 2005). However, it is worth  
 636 pointing out that analysts must confirm the findings through appropriate statistical and  
 637 scientific procedures. In this regard, the insight obtained through the user study with  
 638 chronic conditions data only provides an initial guidance for a further analysis, which is  
 639 clearly out of the scope of the paper.

## 640 6. Conclusions

641 This paper has introduced and analyzed a multivariate visualization method called  
 642 SRA, which is based on a set of radial axis vectors that represent data features, and can gen-

643 erate any linear projection of high-dimensional data points onto a two-dimensional plane.  
644 On the one hand, unlike SC, SRA plots allow users to approximate high-dimensional data  
645 values. On the other hand, in comparison with ARA, SRA provides less cluttered plots,  
646 and allows users to analyze the axis vectors and all of the projected points simultaneously.  
647 Moreover, in SRA longer axis vectors generally represent features that have a smaller in-  
648 fluence on a projection. Since it is easier to identify these vectors, the technique can be  
649 used to carry out an interactive backwards feature selection effectively, where users pro-  
650 gressively eliminate vectors from the plots. Additionally, in contrast to other works in the  
651 literature, we argue that analysts should consider not only the lengths of the axis vectors,  
652 but also their orientations, and expert domain knowledge.

653 In particular, we have used SRA to carry out visual feature selection procedures with  
654 a real-world data set associated with medical chronic conditions of high prevalence in our  
655 society. Results show that SRA allows us to visualize groups of chronic patients with one  
656 or two chronic conditions (DM and/or HBP), while showing the contribution of different  
657 clinical features for discriminating among health statuses. These kinds of visualizations,  
658 which in principle are designed for performing exploratory data analyses, can be very  
659 valuable for experts in the clinical domain. In particular, the visual identification of drugs  
660 and diagnoses somehow related to chronic conditions may be of great value for a better  
661 understanding of these conditions, and may even reveal potential new relationships among  
662 diagnoses and drugs. Therefore, the method proposed in this work can be of great help  
663 to clinicians and health managers for planning care and health resources allocation. This  
664 could lead to an improvement of the health care system, both from an economical and  
665 social point of view.

666 Finally, as future research, we plan to work with time series data in order to find chronic  
667 patient trajectories. This could allow experts to identify the risk factors associated with the  
668 onset or evolution of a chronic condition. As a consequence, health managers could estab-  
669 lish prevention programs according to the risk of a patient of suffering certain conditions.

## 670 **Acknowledgements**

671 This work has been partly funded by the Spanish Ministry of Economy (projects  
672 TIN2014-62143-EXP, TIN2015-70799-R, TEC2016-75361-R, TEC2016-75161-C2-1-4,  
673 and TIN2015-66731-C2-1-R) and the Institute of Health Carlos III (grant DTS17/00158).

## 674 **References**

675 Alcalá-Fdez, J., Sánchez, L., García, S., del Jesús, M. J., Ventura, S., Garrell, J. M.,  
676 Otero, J., Romero, C., Bacardit, J., Rivas, V. M., Fernández, J. C., & Herrera, F. (2008).

- 677 Keel: A software tool to assess evolutionary algorithms for data mining problems. *Soft*  
678 *Comput.*, 13, 307–318.
- 679 Amar, R., Eagan, J., & Stasko, J. (2005). Low-level components of analytic activity in  
680 information visualization. In *Proceedings of the Proceedings of the 2005 IEEE Sympo-*  
681 *sium on Information Visualization* (pp. 15–21). Washington, DC, USA: IEEE Computer  
682 Society.
- 683 Averill, R. F., Goldfield, N., Eisenhandler, J., Muldoon, J. H., Hughes, J., Neff, J. M.,  
684 Gay, J. C., Gregg, L. W., Gannon, D., Shafir, B., Bagadia, F., & Steinbeck, B. (1999).  
685 Development and evaluation of clinical risk groups (crgs). *Wallingford, CT: 3M Health*  
686 *Information Systems*, .
- 687 Berlinguet, M., Preyra, C., & Dean, S. (2005). *Comparing the Value of Three Main*  
688 *Diagnostic-Based Risk-Adjustment Systems (DBRAS)*. Canadian Health Services Re-  
689 search Foundation.
- 690 Bertini, E., Tatu, A., & Keim, D. (2011). Quality metrics in high-dimensional data visu-  
691 alization: An overview and systematization. *IEEE Transactions on Visualization and*  
692 *Computer Graphics*, 17, 2203–2212.
- 693 Blum, A. L., & Langley, P. (1997). Selection of relevant features and examples in machine  
694 learning. *Artificial Intelligence*, 97, 245–271.
- 695 Centers for Disease Control and Prevention (2011). International Classifica-  
696 tion of Diseases, Ninth Revision, Clinical Modification (ICD-9-CM). [Online]  
697 <http://www.cdc.gov/nchs/icd/icd9cm.htm>. Accessed Jan. 2018.
- 698 Chandrashekar, G., & Sahin, F. (2014). A survey on feature selection methods. *Computers*  
699 *and Electrical Engineering*, 40, 16–28.
- 700 Chen, K., & Liu, L. (2004). VISTA: validating and refining clusters via visualization.  
701 *Information Visualization*, 3, 257–270.
- 702 Cox, T., & Cox, M. (1994). *Multidimensional Scaling*. Monographs on Statistics and  
703 Applied Probability 88. Chapman & Hall.
- 704 Dasarathy, B. V. (1991). *Nearest Neighbor (NN) Norms: NN Pattern Classification Tech-*  
705 *niques*. Los Alamitos, CA: IEEE Computer Society Press.
- 706 Ding, C., & Li, T. (2007). Adaptive dimension reduction using discriminant analysis and  
707 k-means clustering. In *Proceedings of the 24th International Conference on Machine*  
708 *Learning ICML'07* (pp. 521–528). New York, NY, USA: ACM.

- 709 Draper, G. M., Livnat, Y., & Riesenfeld, R. F. (2009). A survey of radial methods for  
710 information visualization. *IEEE Transactions on Visualization and Computer Graphics*,  
711 *15*, 759–776.
- 712 Duda, R. O., Hart, P. E., & Stork, D. G. (2001). *Pattern Classification*. Wiley.
- 713 Fernández-Sánchez, J., Soguero-Ruiz, C., de Miguel-Bohoyo, P., Rivas-Flores, F. J., Ángel  
714 Gómez-Delgado, Gutiérrez-Expósito, F. J., & Mora-Jiménez, I. (2017). Clinical risk  
715 groups analysis for chronic hypertensive patients in terms of icd9-cm diagnosis codes. In  
716 *Proceedings of the 4th International Conference on Physiological Computing Systems -*  
717 *Volume 1: PhyCS* (pp. 13–22). INSTICC SciTePress.
- 718 Gabriel, K. R. (1971). The biplot graphic display of matrices with application to principal  
719 component analysis. *Biometrika*, *58*, 453–467.
- 720 Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine*  
721 *Learning*, *63*, 3–42.
- 722 Goldberger, J., Roweis, S., Hinton, G., & Salakhutdinov, R. (2005). Neighborhood compo-  
723 nent analysis. In *Advances in Neural Information Processing Systems 17* (pp. 513–520).
- 724 Gower, J., Gardner-Lubbe, S., & le Roux, N. (2011). *Understanding Biplots*. John Wiley  
725 & Sons.
- 726 Guo, D. (2003). Coordinating computational and visual approaches for interactive feature  
727 selection and multivariate clustering. *Information Visualization*, *2*, 232–246.
- 728 Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Jour-*  
729 *nal of Machine Learning Research*, *3*, 1157–1182.
- 730 Guyon, I., Weston, J., Barnhill, S., & Vapnik, V. (2002). Gene selection for cancer classi-  
731 fication using support vector machines. *Machine Learning*, *46*, 389–422.
- 732 Hughes, J. S., Averill, R. F., Eisenhandler, J., Goldfield, N. I., Muldoon, J., Neff, J. M.,  
733 & Gay, J. C. (2004). Clinical Risk Groups (CRGs): a classification system for risk-  
734 adjusted capitation-based payment and health care management. *Medical care*, *42*, 81–  
735 90.
- 736 Hyvärinen, A., Karhunen, J., & Oja, E. (2001). *Independent component analysis*. Adaptive  
737 and learning systems for signal processing, communications, and control. J. Wiley.

- 738 Ingram, S., Munzner, T., Irvine, V., Tory, M., Bergner, S., & Möller, T. (2010). Dimstiller:  
739 Workflows for dimensional analysis and reduction. In *IEEE VAST* (pp. 3–10). IEEE  
740 Computer Society.
- 741 Inselberg, A., & Dimsdale, B. (1990). Parallel coordinates: a tool for visualizing multi-  
742 dimensional geometry. In *Proceedings of the 1st conference on Visualization VIS'90*  
743 (pp. 361–378). Los Alamitos, CA, USA: IEEE Computer Society Press.
- 744 Johansson, S., & Johansson, J. (2009). Interactive dimensionality reduction through user-  
745 defined combinations of quality metrics. *IEEE Transactions on Visualization & Com-  
746 puter Graphics, 15*, 993–1000.
- 747 Jolliffe, I. T. (2010). *Principal component analysis*. Springer series in statistics. Springer-  
748 Verlag.
- 749 Kandogan, E. (2000). Star coordinates: A multi-dimensional visualization technique with  
750 uniform treatment of dimensions. In *In Proceedings of the IEEE Information Visualiza-  
751 tion Symposium, Late Breaking Hot Topics* (pp. 9–12).
- 752 Kandogan, E. (2001). Visualizing multi-dimensional clusters, trends, and outliers using  
753 star coordinates. In *Proceedings of the seventh ACM SIGKDD international conference  
754 on Knowledge discovery and data mining KDD'01* (pp. 107–116). New York, NY, USA:  
755 ACM.
- 756 Krause, J., Perer, A., & Bertini, E. (2014). Infuse: Interactive feature selection for predic-  
757 tive modeling of high dimensional data. *IEEE Transactions on Visualization & Com-  
758 puter Graphics, 20*, 1614–1623.
- 759 Leban, G., Zupan, B., Vidmar, G., & Bratko, I. (2006). VizRank: Data Visualization  
760 Guided by Machine Learning. *Data Mining and Knowledge Discovery, 13*, 119–136.
- 761 Lichman, M. (2013). UCI machine learning repository.
- 762 Maaten, L. v. (2015). Matlab toolbox for dimensionality reduction.
- 763 MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate  
764 observations. In L. L. Cam, & J. Neyman (Eds.), *Proc. of the fifth Berkeley Symposium  
765 on Mathematical Statistics and Probability* (pp. 281–297). University of California  
766 Press volume 1.
- 767 May, T., Bannach, A., Davey, J., Ruppert, T., & Kohlhammer, J. (2011). Guiding feature  
768 subset selection with an interactive visualization. In *2011 IEEE Conference on Visual  
769 Analytics Science and Technology (VAST)* (pp. 111–120).

- 770 McLachlan, G. J. (2004). *Discriminant analysis and statistical pattern recognition*. Wiley  
771 series in probability and mathematical statistics. Probability and mathematical statistics.  
772 Wiley-Interscience.
- 773 McNicoll, G. (2002). World population ageing 1950-2050. *Population and Development*  
774 *Review*, 28, 814–816.
- 775 Norwegian Institute of Public Health (2017). *WHO Collaborating Centre for Drug Statis-*  
776 *tics Methodology, Guidelines for ATC classification and DDD assignment 2018*. Oslo.
- 777 Organization, W. H. (1999). Hypertension guidelines. *J Hypertension*, 17, 151–183.
- 778 Organization, W. H. et al. (2005). *Preventing Chronic Diseases-A Vital Investment: WHO*  
779 *Global Report*. World Health Organization.
- 780 Paulovich, F. V., Nonato, L. G., Minghim, R., & Levkowitz, H. (2008). Least square pro-  
781 jection: A fast high-precision multidimensional projection technique and its application  
782 to document mapping. *IEEE Transactions on Visualization and Computer Graphics*,  
783 14, 564–575.
- 784 Rauber, P. E., Silva, R. R. O. d., Feringa, S., Celebi, M. E., Falcão, A. X., & Telea,  
785 A. C. (2015). Interactive Image Feature Selection Aided by Dimensionality Reduction.  
786 In *EuroVis Workshop on Visual Analytics (EuroVA)* (pp. 67–74). The Eurographics  
787 Association.
- 788 Rauber, T., & Steiger-Garçon, A. (1993). Feature selection of categorical attributes based  
789 on contingency table analysis. In *Proceedings of the 5th Portuguese Conference on*  
790 *Pattern Recognition, Porto, Portugal*.
- 791 Rubio-Sánchez, M., Raya, L., Díaz, F., & Sanchez, A. (2016). A comparative study be-  
792 tween radviz and star coordinates. *IEEE Transactions on Visualization and Computer*  
793 *Graphics*, 22, 619–628.
- 794 Rubio-Sánchez, M., & Sanchez, A. (2014). Axis calibration for improving data attribute  
795 estimation in star coordinates plots. *IEEE Transactions on Visualization and Computer*  
796 *Graphics*, 20, 2013–2022.
- 797 Rubio-Sánchez, M., Sanchez, A., & Lehmann, D. J. (2017). Adaptable radial axes plots for  
798 improved multivariate data visualization. *Computer Graphics Forum (Proc. EuroVis)*, .
- 799 Seo, J., & Shneiderman, B. (2005). A rank-by-feature framework for interactive explo-  
800 ration of multidimensional data. *Information Visualization*, 4, 96–113.



- 801 Shneiderman, B. (1996). The eyes have it: A task by data type taxonomy for information  
802 visualizations. In *Proceedings of the 1996 IEEE Symposium on Visual Languages* (pp.  
803 336–343). Washington, DC, USA: IEEE Computer Society.
- 804 Soguero-Ruiz, C., Hindberg, K., Mora-Jiménez, I., Rojo-Álvarez, J. L., Skrøvseth, S. O.,  
805 Godtlielsen, F., Mortensen, K., Revhaug, A., Lindsetmo, R.-O., Augestad, K. M. et al.  
806 (2016). Predicting colorectal surgical complications using heterogeneous clinical data  
807 and kernel methods. *Journal of biomedical informatics*, *61*, 87–96.
- 808 Sun, Y., Yuan, J., Hu, Y., & Xiao, W. (2008). An improved multivariate data visualization  
809 technique. In *International Conference on Information and Automation, ICIA'08*. (pp.  
810 1525–1530).
- 811 Tatu, A., Maaß, F., Färber, I., Bertini, E., Schreck, T., Seidl, T., & Keim, D. A. (2012).  
812 Subspace search and visualization to make sense of alternative clusterings in high-  
813 dimensional data. In *Proc. IEEE Symposium on Visual Analytics Science and Tech-*  
814 *nology* (pp. 63–72). IEEE Computer Society.
- 815 Tsai, C.-Y., & Chiu, C.-C. (2008). A clustering-oriented star coordinate translation method  
816 for reliable clustering parameterization. In *Proceedings of the 12th Pacific-Asia confer-*  
817 *ence on Advances in knowledge discovery and data mining PAKDD'08* (pp. 749–758).  
818 Berlin, Heidelberg: Springer-Verlag.
- 819 Wang, Y., Li, J., Nie, F., Theisel, H., Gong, M., & Lehmann, D. J. (2017). Linear discrim-  
820 inative star coordinates for exploring class and cluster separation of high dimensional  
821 data. *Computer Graphics Forum (Proc. EuroVis)*, .
- 822 Weinberger, K. Q., & Saul, L. K. (2009). Distance metric learning for large margin nearest  
823 neighbor classification. *Journal of Machine Learning Research*, *10*, 207–244.
- 824 Witten, I. H., & Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and*  
825 *Techniques*. (2nd ed.). San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- 826 Yang, J., Peng, W., Ward, M. O., & Rundensteiner, E. A. (2003a). Interactive hierarchical  
827 dimension ordering, spacing and filtering for exploration of high dimensional datasets.  
828 In *Proceedings of the Ninth Annual IEEE Conference on Information Visualization IN-*  
829 *FOVIS'03* (pp. 105–112). Washington, DC, USA: IEEE Computer Society.
- 830 Yang, J., Ward, M. O., & Rundensteiner, E. A. (2002). Interring: An interactive tool  
831 for visually navigating and manipulating hierarchical structures. In *Proceedings of the*  
832 *IEEE Symposium on Information Visualization (InfoVis'02) INFOVIS'02* (pp. 77–84).  
833 IEEE Computer Society.

- 834 Yang, J., Ward, M. O., & Rundensteiner, E. A. (2003b). Interactive hierarchical displays:  
835 A general framework for visualization and exploration of large multivariate data sets.  
836 *Computers & Graphics*, (pp. 265–283).
- 837 Yi, J. S., ah Kang, Y., Stasko, J., & Jacko, J. (2007). Toward a deeper understanding of  
838 the role of interaction in information visualization. *IEEE Transactions on Visualization*  
839 *and Computer Graphics*, 13, 1224–1231.
- 840 Zhang, K.-B., Orgun, M., & Zhang, K. (2006). HOV<sup>3</sup>: An approach to visual cluster  
841 analysis. In *Advanced Data Mining and Applications* (pp. 316–327). Springer Berlin /  
842 Heidelberg volume 4093 of *Lecture Notes in Computer Science*.