

# Accepted Manuscript

Improved multiclass feature selection via list combination

Javier Izetta, Pablo F. Verdes, Pablo M. Granitto

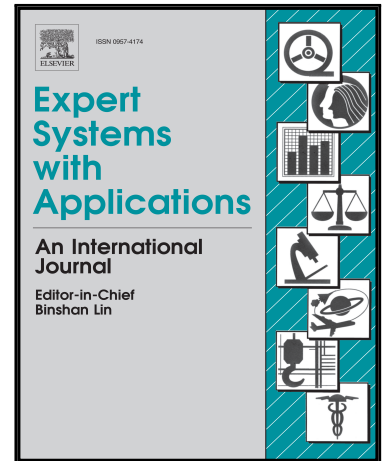
PII: S0957-4174(17)30467-0  
DOI: [10.1016/j.eswa.2017.06.043](https://doi.org/10.1016/j.eswa.2017.06.043)  
Reference: ESWA 11414

To appear in: *Expert Systems With Applications*

Received date: 8 March 2017  
Revised date: 30 June 2017  
Accepted date: 30 June 2017

Please cite this article as: Javier Izetta, Pablo F. Verdes, Pablo M. Granitto, Improved multi-class feature selection via list combination, *Expert Systems With Applications* (2017), doi: [10.1016/j.eswa.2017.06.043](https://doi.org/10.1016/j.eswa.2017.06.043)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



**Highlights**

- We introduce new SVM-RFE feature selection methods for multiclass problems
- We use binary decomposition followed by strategies to combine lists of features
- We discuss statistical approaches and voting theory methods
- One-vs-One methods give better results than One-vs-All methods
- The new K-First method is the more effective in selecting relevant features

# Improved multiclass feature selection via list combination

Javier Izetta<sup>a</sup>, Pablo F. Verdes<sup>a</sup>, Pablo M. Granitto<sup>a,\*</sup>

<sup>a</sup>*CIFASIS, French Argentine International Center for Information and Systems  
Sciences, UNR-CONICET,  
Bv. 27 de Febrero 210 Bis, 2000 Rosario, Argentina*

---

## Abstract

Feature selection is a crucial machine learning technique aimed at reducing the dimensionality of the input space. By discarding useless or redundant variables, not only it improves model performance but also facilitates its interpretability. The well-known Support Vector Machines-Recursive Feature Elimination (SVM-RFE) algorithm provides good performance with moderate computational efforts, in particular for wide datasets. When using SVM-RFE on a multiclass classification problem, the usual strategy is to decompose it into a series of binary ones, and to generate an importance statistics for each feature on each binary problem. These importances are then averaged over the set of binary problems to synthesize a single value for feature ranking. In some cases, however, this procedure can lead to poor selection. In this paper we discuss six new strategies, based on list combination, designed to yield improved selections starting from the importances given by the binary problems. We evaluate them on artificial and real-world datasets, using both One-Vs-One (OVO) and One-Vs-All (OVA) strategies. Our results suggest that the OVO decomposition is most effective for feature selection on multiclass problems. We also find that in most situations the new K-First strategy can find better subsets of features than the traditional weight average approach.

*Keywords:*

---

\*Corresponding author

*Email addresses:* [izetta@cifasis-conicet.gov.ar](mailto:izetta@cifasis-conicet.gov.ar) (Javier Izetta),  
[verdes@cifasis-conicet.gov.ar](mailto:verdes@cifasis-conicet.gov.ar) (Pablo F. Verdes),  
[granitto@cifasis-conicet.gov.ar](mailto:granitto@cifasis-conicet.gov.ar) (Pablo M. Granitto)

Feature Selection, Multiclass problems, Support Vector Machine

---

ACCEPTED MANUSCRIPT

## 1 1. Introduction

2 Many important problems in Machine Learning, as well as in in-silico  
3 Chemistry (Raies & Bajic, 2016), Biology, “high-throughput” technologies  
4 (Golub et al., 1999; Leek et al., 2010) or text processing (Forman, 2003;  
5 Uysal, 2016), share the property of involving much more features than mea-  
6 sured samples are available (Guyon & Elisseeff, 2003). The datasets associ-  
7 ated to these problems are, unsurprisingly, called “wide”. Usually, most of  
8 these variables carry a relatively low importance for the problem at hand.  
9 Furthermore, in some cases they interfere with the learning process instead  
10 of helping it, a scenario usually referred to as “curse of dimensionality”.

11 Feature selection is an important pre-processing technique of Machine  
12 Learning aimed at coping with this curse (Kohavi & John, 1997). Its main  
13 goal is to find a small subset of the measured variables that improve, or at  
14 least do not degrade, the performance of the modeling method applied to the  
15 dataset. But feature selection methods do not only avoid the curse of dimen-  
16 sionality: they also allow for a considerable reduction in model complexity,  
17 an easier visualization and, in particular, a better interpretation of the data  
18 under analysis and the developed models (Liu et al., 2005).

19 Several methods have been introduced in recent years, from general ones  
20 like Wrappers (Kohavi & John, 1997) and filters (Kira & Rendell, 1992) to  
21 very specific ones developed for SVM (Weston et al., 2000; Nguyen & De la  
22 Torre, 2010) and RVM (Mohsenzadeh et al., 2013, 2016) classifiers. Amongst  
23 other methods in the field (Hua et al., 2009), the well-known Recursive Fea-  
24 ture Elimination (RFE) algorithm provides good performance with moderate  
25 computational efforts (Guyon et al., 2002) on wide datasets. The original and  
26 most popular version of this method uses a linear Support Vector Machine  
27 (SVM) (Vapnik, 2013) to select the candidate features to be eliminated. Ac-  
28 cording to the SVM-RFE algorithm, the importance of an input variable  
29  $i$  is directly correlated with the corresponding component ( $w_i$ ) of the vec-  
30 tor defining the separating hyperplane ( $\mathbf{w}$ ). The method is widely used in  
31 Bioinformatics (Guyon et al., 2002; Statnikov et al., 2005). Alternative RFE  
32 methods using other classifiers have also been introduced in the literature  
33 (Granitto et al., 2006; You et al., 2014).

34 Typical feature selection algorithms are designed for binary classification  
35 problems, as the original version of RFE. Multiclass problems have received  
36 much less attention because of their increased difficulty. Also, because some  
37 classifiers involved in the selection process are designed to solve binary prob-

38 lems. Most methods available for feature selection on multiclass problems  
 39 are simple extensions of base methods. For example, RFE can be associated  
 40 to a multiclass classifier like Random Forest (Breiman, 2001; Granitto et al.,  
 41 2006).

42 Although SVM was originally developed to deal only with binary prob-  
 43 lems, it was extended to directly solve multiclass problems in different man-  
 44 ners (Weston & Watkins, 1999; Crammer & Singer, 2001; Hsu & Lin, 2002),  
 45 but with a modest success attributed mainly to the increased complexity of  
 46 the solutions. On the other hand, in the last years several methods were  
 47 developed to solve a multiclass problem using an appropriate combination of  
 48 binary classifiers (Allwein et al., 2000; Hsu & Lin, 2002). The most usually  
 49 followed strategy for multiclass SVM is known as “One-vs-One” (OVO). Ac-  
 50 cording to this approach, a classification problem with  $c$  classes is replaced  
 51 with  $M = c(c - 1)/2$  reduced binary ones, each one of them consisting of dis-  
 52 criminating a pair of classes. In order to classify a new example, it is passed  
 53 through all binary classifiers and the most voted class is selected. Another  
 54 useful strategy is “One-vs-All” (OVA). In this second case, a problem with  
 55  $c$  classes is replaced with  $M = c$  reduced binary problems, each one of them  
 56 consisting of discriminating a single class from all remaining ones.

57 Therefore, the most usual approach to implement a multiclass SVM-RFE  
 58 method is to directly apply the RFE algorithm over an OVO or OVA multi-  
 59 class SVM (Ramaswamy et al., 2001; Duan et al., 2007; Zhou & Tuck, 2007).  
 60 The pioneering work of Ramaswamy et al. (Ramaswamy et al., 2001) pro-  
 61 posed the OVA solution, but also compared results with the OVO strategy.  
 62 Duan et al. (Duan et al., 2007) and Zhou et al. (Zhou & Tuck, 2007) devel-  
 63 oped slight variations of the method, always considering both OVA and OVO  
 64 implementations. Zhou et al. (Zhou & Tuck, 2007) also considered solutions  
 65 to the RFE problem using a direct multiclass implementation.

66 Interestingly, the solutions to the multiclass SVM-RFE problem that we  
 67 have just described involve an important decision about the feature selection  
 68 process which is usually neglected: they rank features by simply averag-  
 69 ing components over the binary problems. For an input variable  $i$  they use  
 70  $\langle |w_{ij}| \rangle_j$ , the mean importance over all binary problems  $j$ , as the corre-  
 71 sponding importance. As we discuss in the next section, this strategy can  
 72 lead to sub-optimal selections in many cases. Once the original multiclass  
 73 problem has been divided into multiple binary ones, the feature selection  
 74 problem can be treated in a similar way. Then, a possible solution is to cast  
 75 the multiclass feature selection problem as the problem of selecting candidate

76 features from multiple lists (Jurman et al., 2008), each list corresponding to  
77 a different binary sub-problem.

78 Similar solutions have been studied in related fields. In Bioinformatics,  
79 for example, Haury et al. (Haury et al., 2011) discussed the combination of  
80 multiple lists of genes from bootstraps of the same gene-expression dataset.  
81 Zhou and Dickerson (Zhou & Dickerson, 2014) and Zhou and Wang (Zhou &  
82 Wang, 2016) proposed the use of class-dependent features (different features  
83 for each binary problem) for biomarker discovery. Dittman et al. (Dittman  
84 et al., 2013) showed that combining multiple lists in binary classification  
85 problems can improve the feature selection results. In a short work in text  
86 categorization, Neumayer et al. (Neumayer et al., 2011) suggested that the  
87 combination of rankings generated by diverse methods can improve the re-  
88 sults of using a single method. Kanth and Saraswathi (Kanth & Saraswathi,  
89 2015) used class-dependent features for speech emotion recognition, but us-  
90 ing independent features for each class, not a final unique list.

91 In this work we discuss in depth the use of combination of multiple lists in  
92 feature selection for multiclass classification problems. We first introduce a  
93 simple mathematical framework for multiple lists. Using this framework, we  
94 propose diverse strategies to produce improved selection of feature subsets  
95 with SVM-RFE. Also, we use some specifically-designed artificial datasets  
96 and real-world examples to evaluate them extensively, using both the OVO  
97 and OVA strategies.

98 The rest of this article is organized as follows: in Section 2, we describe the  
99 feature selection methods introduced in this work. In Section 3 we evaluate  
100 these methods on diverse datasets and experimental setups. Finally, we draw  
101 our conclusions in Section 4.

## 102 **2. List combination methods for SVM-RFE**

103 The RFE selection method is a recursive process that ranks variables  
104 according to a given importance measure. At each iteration of the algo-  
105 rithm, the importance of each feature is calculated and the less relevant one  
106 is removed—in order to speed up the process, not one but a group of low  
107 relevance features is usually removed. Recursion is needed because the rel-  
108 ative importance of each feature can change substantially when evaluated  
109 over a different subset of features during the stepwise elimination process, in  
110 particular for highly correlated features. The inverse order in which features

111 are eliminated is used to create a final ranking. Then, the feature selection  
 112 process itself is reduced to take the first  $n$  features from this ranking.

113 In the original binary version of SVM-RFE (Guyon et al., 2002), the  
 114 projection of  $\mathbf{w}$  (the normal vector to SVM’s decision hyperplane) in the  
 115 direction of feature  $i$ ,  $w_i$ , is used as the importance measure. The method  
 116 was efficiently extended to multiclass problems, employing the well-known  
 117 OVO or OVA strategies to decompose the multiclass problem into a series  
 118 of related binary ones (Ramaswamy et al., 2001; Duan et al., 2007; Zhou &  
 119 Tuck, 2007). In both cases a set of  $M$  related binary problems is generated,  
 120 each one solved by a vector  $\mathbf{w}_j$ . For each binary problem  $j$ , the importance  
 121 of feature  $i$  is given by the corresponding component,  $w_{ij}$ .

122 In order to obtain a unique importance for each feature in this setup, the  
 123 simplest solution is to average the absolute value of the components  $|w_{ij}|$   
 124 over all related binary problems. We will call this method “Average” in the  
 125 following. The Average solution is implemented, to the best of our knowledge,  
 126 in all available RFE software packages, including the most popular amongst  
 127 researchers (MATLAB, R and PYTHON platforms).

128 However, the only real advantage of the Average strategy is its simplicity.  
 129 Two main drawbacks of this approach should be taken into consideration but  
 130 are usually ignored:

- 131 1. The first issue can be called the *flattening* problem. Consider, for  
 132 example, a feature  $e$  which is able to separate class  $j$  from all remaining  
 133 classes, but is uninformative in other cases. Component  $w_{ej}$  will be  
 134 large, but components  $w_{ek}$  with  $k \neq j$  will be small, giving a low value  
 135 for  $\langle w_{ej} \rangle_j$ . Consider now another feature  $d$  which can give a modest  
 136 help in separating any class from the others, obtaining always moderate  
 137 values of  $w_{dj}$ , and therefore giving a medium value for  $\langle w_{dj} \rangle_j$ . The  
 138 Average strategy will clearly rank the latter over the former, but in  
 139 most scenarios it will be desirable to keep the first variable over the  
 140 second.
- 141 2. The second issue with the Average solution refers to *relative scales*.  
 142 The length of vector  $\mathbf{w}_j$  is different for each binary problem, as it de-  
 143 pends on the margin of the solution, which can change considerably for  
 144 classes that are relatively close or far away in feature space. Averaging  
 145 components of vectors of different lengths can lead to the selection of  
 146 sub-optimal subsets.

New strategies for feature selection able to overcome these drawbacks are



Ranking	List 1	List 2	...	List M
1	$f_2$	$f_3$	...	$f_1$
2	$f_1$	$f_7$	...	$f_3$
3	$f_5$	$f_2$	...	$f_6$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$p$	$f_8$	$f_4$	...	$f_7$

Table 1: List of ranked features for each binary problem.

needed. Here we propose to cast the problem as a selection of candidate features from multiple ranking lists (Jurman et al., 2008). We start by decomposing the multiclass problem into a set of  $M$  related binary problems (through the OVA or OVO strategies). The problem involves a set of  $p$  features,  $F = \{f_1, f_2, \dots, f_p\}$ . SVM-RFE produces a ranking (an ordered list) for each individual problem using the components  $w_{ij}$ . An example is shown in Table 1. This set of lists can be arranged in a matrix (Table 2) where each row shows the position of each feature in the ranking produced for the binary problem shown on each column. We can now define a matrix of relative ranking positions as:

$$r_{i,j} = \frac{p - pos_{i,j}}{p} = 1 - \frac{pos_{i,j}}{p},$$

147 where  $r_{i,j}$  is the relative ranking of feature  $f_i$  in the list corresponding to bi-  
 148 nary problem  $j$ ,  $pos_{i,j}$  is the position of the same feature in the corresponding  
 149 ranking (Table 2), and  $p$  is the total number of features in the problem. No-  
 150 tice that the values of  $r_{i,j}$  belong to the unit interval  $[0, 1]$  and depend linearly  
 151 on the ranking position (a value of  $1 - 1/p$  must be interpreted as the first  
 152 position in the ranking).

153 Important features should reflect in high values of  $r_{i,j}$  for some  $i, j$ , mean-  
 154 ing they are relevant to at least some of the binary classification problems.  
 155 Two main strategies can be used to select those relevant features from this  
 156 matrix and are discussed in the following.

### 157 2.1. Methods based on relative ranking statistics

158 The first strategy consists of measuring an appropriate statistic for each  
 159 feature over all binary problems, and then using it to elaborate a final ranking  
 160 of features. We selected the following four methods:

Features	List 1	List 2	...	List M
$f_1$	2	5	...	1
$f_2$	1	4	...	7
$f_3$	4	1	...	2
$\vdots$	$\vdots$	$pos_{i,j}$	$\ddots$	$\vdots$
$f_p$	$pos_{p,1}$	...	...	$pos_{p,M}$

Table 2: Matrix showing the position of each feature in the ranking of each binary problem. Rows correspond to features and columns to binary problems.

### 161 2.1.1. Average-SD

162 In this method, feature ranking is given by the average value of the relative  
163 position over all binary problems:

$$R_i = \langle r_{i,j} \rangle_j,$$

164 where  $R_i$  is the ranking of feature  $f_i$  in the final ordered list, used to select  
165 features in the multiclass problem. Ties are broken by the standard deviation  
166 (SD) of the relative position (higher is better). We show in the next section  
167 that features with higher SD are preferable over lower SD ones, because a  
168 larger SD means that the feature has some better-than-average rankings.

169 Average-SD can be considered as the base strategy for multiple lists. It  
170 can overcome the *relative scales* problem on averaging weights, but is not  
171 expected to solve the *flattening* problem.

### 172 2.1.2. Best Ranking

173 In this second approach we rank every feature according to the best rel-  
174 ative ranking that it reaches over the set of binary problems:

$$R_i = \max(r_{i,j})_j.$$

175 Ties are broken by the mean value of the relative position over all prob-  
176 lems. A similar method has been used to select the winning class in multiple  
177 classifier systems (Ho et al., 1994). This strategy can be viewed as an ex-  
178 treme case, considering for each feature just one of the multiple rankings it  
179 receives and disregarding the rest. On the other hand, it is most aggressive  
180 in dealing with the flattening problem.

181 *2.1.3. 3Q-SD*

182 The third method orders features according to the 3rd quartile of the  
183 distribution of relative rankings:

$$R_i = 3Q(r_{i,j})_j,$$

184 where the 3Q function returns the 3rd quartile of its argument. As in  
185 Average-SD, ties are broken by the SD. This approach is intermediate be-  
186 tween the two previous ones, searching for features that reach a high relative  
187 position, but also considering the full relative rankings distribution.

188 *2.1.4. K-First*

189 This method is adapted from a strategy to select relevant documents in  
190 information retrieval (Nuray & Can, 2006). The idea is to only consider  
191 features located in the top  $k$  positions of each individual list. We re-scale the  
192 relative rankings with a linear mapping reaching 0 for the  $k + 1$  feature, and  
193 then take the average of this new relative importance:

$$r'_{i,j} = \max\left(1 - \frac{pos_{i,j}}{k}, 0\right)$$

$$R_i = \langle r'_{i,j} \rangle_j,$$

194 where  $r'_{i,j}$  is the re-scaled relative weight for feature  $f_i$  and  $k$  is the number  
195 of features to be considered from each list ( $k < p$ ). As in the Best Ranking  
196 method, ties are broken by the mean value of the original relative ranking,  
197  $\langle r_{i,j} \rangle_j$ . We discuss the set of parameter values  $k$  in the next section.  
198 This strategy is aimed at searching for features which are highly relevant for  
199 some of the problems, but is not limited to searching for the most relevant  
200 features —as the Best Ranking method is. It can potentially overcome both  
201 drawbacks of Average: relative scaling and flattening.

202 *2.2. Methods based on voting theory*

203 The second general strategy is related to voting theory (Saari, 2001;  
204 Young, 1988). In this setup we consider each binary problem as a voter,  
205 producing a ranking over a set of  $p$  candidates. Multiple methods were de-  
206 veloped over the years to solve the problem of combining elector preferences  
207 to find winner candidates —the most useful of them are known collectively  
208 as "Condorcet Methods". We focused on two popular procedures as selection  
209 methods for relevant features over multiple lists:

210 *2.2.1. Condorcet*

211 The most basic Condorcet method is known as Copeland’s method, or  
 212 simply as Condorcet method (we will use the latter name in this work). It  
 213 confronts each pair of features on every list (all binary problems), and then  
 214 counts the number of wins minus the number of defeats for each feature  
 215 (Young, 1988). A feature wins over another if it is ranked higher in the  
 216 considered list. The global difference between wins and loses is used to  
 217 rank features in the multiclass problem. Ties are broken by average relative  
 218 rankings.

219 *2.2.2. Schulze*

220 This method, introduced by Schulze informally in 1997 and published  
 221 later (Schulze, 2011), represents an improvement over previous Condorcet  
 222 methods. It begins by counting wins and loses over each pair of features and  
 223 all lists, storing these numbers in a pairwise preference matrix. Then a graph  
 224 is constructed, with features as nodes and values in the matrix as weights.  
 225 Finally, using a variant of the FloydWarshall algorithm, the strongest paths  
 226 over the graph are selected for each pair of features, and their strengths  
 227 are used to compare features. The strength of a path is defined as that  
 228 of its weakest link (i.e., lower value in the matrix of preferences). A path  
 229 between two nodes is valid if there is a sequence of strictly decreasing weights  
 230 connecting them (Schulze, 2011). Features with more wins upon strength  
 231 comparison are ranked first. The method is expected to perform better than  
 232 basic Condorcet, but the computational load involved is significant.

233 **3. Evaluation on artificial datasets**

234 We first consider artificial classification problems in order to evaluate spe-  
 235 cific aspects of the new methods and to be able to compare their capabilities  
 236 in a controlled manner.

237 *3.1. Experimental setup*

238 As in previous works (Granitto et al., 2006), we strive to use an appro-  
 239 priate computational setup for feature selection. We perform  $n = 100$  times  
 240 a random split (75% – 25%) of each dataset in training and testing sets (the  
 241 former are used to select features and train the classifiers, while the latter  
 242 for model accuracy estimation). The testing sets are completely external to

243 the feature selection process, thus providing unbiased estimates of classifica-  
 244 tion errors for different number of features. The results of the  $n$  replicated  
 245 experiments are then aggregated to yield mean error rate estimations and  
 246 their corresponding SD.

247 SVM-RFE was implemented using the OVO strategy unless specified  
 248 otherwise. In both cases, we created the corresponding binary problems and  
 249 produced a ranking of features for each of them. To create a ranking we used  
 250 the standard SVM-RFE (linear kernel), as described by their authors (Guyon  
 251 et al., 2002), eliminating 10% of the features at each iteration until there were  
 252 less than 20 features left, when we slowed the procedure to eliminate 1 feature  
 253 at each iteration. The fixed set of lists of ranked features were combined  
 254 using the methods described before, producing a final list of features for each  
 255 method under evaluation. Finally, for each method we fitted a multiclass  
 256 SVM for a varying number of features, from 2 to  $p$ , using only the training  
 257 data, and measured the classification error using the testing set. The  $C$   
 258 parameter was estimated in all cases using 5-fold cross validation of the  
 259 training set.

### 260 3.2. Artificial datasets

261 We created three different multiclass datasets that provide diverse chal-  
 262 lenges to our methods. In all cases, each class is sampled from a Gaussian dis-  
 263 tribution with diagonal covariance matrix. For each dataset we can identify,  
 264 by construction, a group of relevant features that can discriminate amongst  
 265 classes and another group of irrelevant features containing Gaussian noise.  
 266 All noisy features have the same mean (0) and SD (1) for all classes. Each  
 267 dataset is composed of 3000 points evenly distributed among classes. The  
 268 number of noisy features is fixed at 500.

269 In the first dataset, called Artificial-1, there is a group of 5 features that  
 270 is relevant for each class, i.e., class-specific features. The set of 5 features  
 271 together shift the class center away from other classes. All relevant features  
 272 have the same importance for the problem. The SD of the Gaussian distri-  
 273 butions corresponding to relevant features are always set to 0.5.

274 A different situation arises when there are sets of features which are rel-  
 275 evant to some of the classes (more than one) but not for all of them. We  
 276 created a second classification problem, Artificial-2, to evaluate this chal-  
 277 lenge. The dataset has 8 classes and 25 relevant features, all sampled from  
 278 Gaussian distributions with a SD of 0.5. The first 5 features are relevant for  
 279 the first 3 classes of the problem. The following 5 are relevant for classes 4

Dataset	Classes	Relevant features	Noisy features
Artificial-1-3C	3	15	500
Artificial-1-4C	4	20	500
Artificial-1-5C	5	25	500
Artificial-1-8C	8	40	500
Artificial-1-16C	16	80	500
Artificial-2-8C	8	25	500
Artificial-3-3C	3	15	500
Artificial-3-8C	8	40	500
Artificial-3-16C	16	80	500

Table 3: Details of the artificial datasets used in this work.

280 and 5 only, and are less relevant than the first 5. The rest of the features are  
 281 relevant for a single class, 5 for each of the remaining 3 classes. These last  
 282 features are less relevant than the first 10 features.

283 Finally, we created a third problem, called Artificial-3, where all relevant  
 284 features are equally useful for all classes at the same time. As in Artificial-1,  
 285 there are 5 features for each class, all sampled from Gaussian distributions  
 286 with a SD of 0.5.

287 In all problems there is an overlap among classes, giving a nonzero Bayes  
 288 error. We created five datasets for Artificial-1, with an increasing number  
 289 of classes, and 3 datasets for Artificial-3 in the same way. Table 3 collects  
 290 technical details of the datasets.

291

### 292 3.3. Methodological setup

#### 293 3.3.1. *K-First*

294 The *K-First* method is the only approach involving a parameter that  
 295 needs to be set,  $k$ . The value of this parameter regulates the number of  
 296 variables that receive a relative ranking. A very low value would make the  
 297 method similar to Best Ranking, while a high one would turn the method  
 298 into 3Q-SD (furthermore,  $k = p$  would convert the method into Average-SD).

299 We evaluated several values of  $k$  (increasing fractions of  $p$ ) over all ar-  
 300 tificial datasets considered. Figure 1 shows the corresponding error curves  
 301 as a function of the number of features selected by the method, for some  
 302 representative problems. Error curves for all other artificial problems are

303 similar to the reported ones (we include more figures in the Additional Ma-  
 304 terial section). The vertical dotted lines show the correct number of relevant  
 305 features for the problem, i.e. where the minimum of the curve should ideally  
 306 be located. When possible, we show with a gray horizontal line the chance  
 307 error level for the corresponding problem. The top row shows typical results  
 308 for the Artificial-3 problem. In this case there are no class-specific features  
 309 and, as a consequence, the results are almost independent of  $k$ . The bottom  
 310 row shows results for class-dependent problems, Artificial-1 and 2. In this  
 311 case the results clearly depend on  $k$ . We found that a value of 10% of  $p$  gives  
 312 consistently good results in all artificial cases considered here, therefore we  
 313 will use this value for the rest of the paper.

### 314 3.3.2. Average-SD and 3Q-SD

315 As noted before, these methods use the SD of the relative rankings as a  
 316 breaking tie criterion, considering larger values of SD as better than smaller  
 317 ones. This is based on the assumption that a large SD is associated with  
 318 high rankings for some of the binary problems, and that such behavior is  
 319 able to highlight class-dependent features over flat ones. In order to confirm  
 320 this, we compared for both methods over a set of artificial problems the  
 321 use of maximum versus minimum SD to break ties. Figure 2 shows the  
 322 corresponding results for some representative cases. They are similar in all  
 323 other cases (some of which are shown in the Additional Material section).  
 324 As this figure shows, using maximum values always leads to equal or better  
 325 performance than using minimum ones.

### 326 3.3.3. OVA-SVM vs. OVO-SVM

327 We applied both OVA and OVO strategies, combined with our feature  
 328 selection methods, to all artificial datasets. We compared all results and  
 329 found that the OVO strategy yields equal or superior performance in all  
 330 cases. In Figure 3 we show some representative examples of this comparison,  
 331 using the *Artificial-1-8C* and *Artificial-2-8C* datasets. On the left column  
 332 we show OVA results, while the OVO case is depicted on the right column.  
 333 We use the same scale for the corresponding panels. We also included the  
 334 Bayes error for both datasets as dotted horizontal lines, and the true number  
 335 of relevant features as a dotted vertical lines. More datasets are included in  
 336 the Additional Material section.

337 It is interesting to note that the two methods more directly aimed at find-  
 338 ing class-relevant features (K-First and Best Ranking) are the ones showing

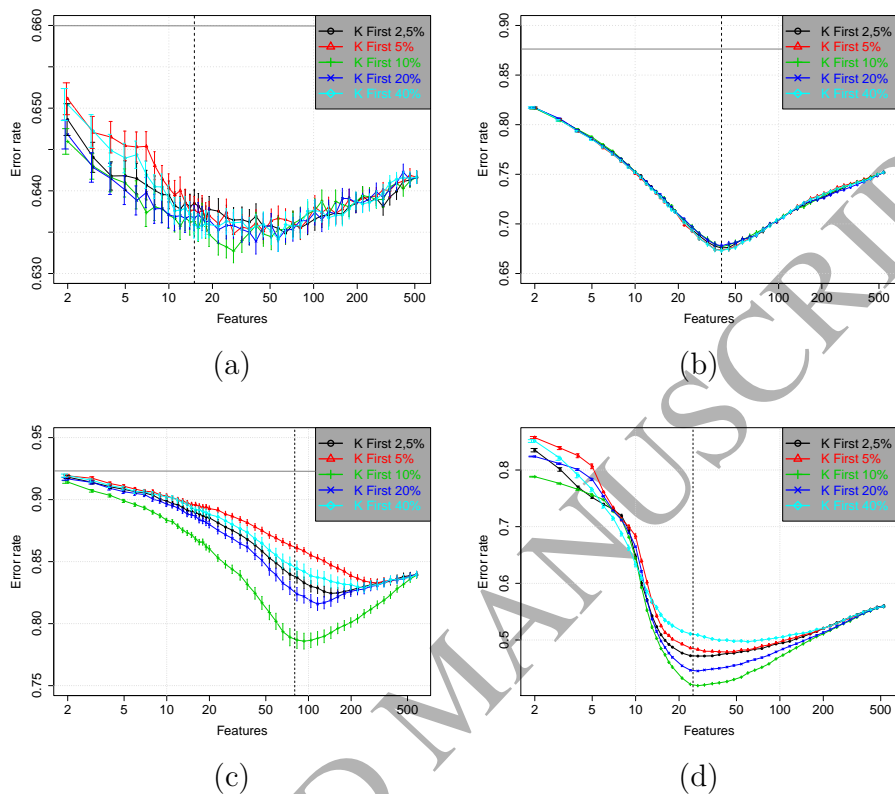


Figure 1: Evaluation of different values of  $k$  for the K-First method. Each line shows average error rates as a function of the number of features selected by the corresponding method, with 1 SD error bars. (a) Artificial-3-3C (b) Artificial-3-8C (c) Artificial-1-16C (d) Artificial-2-8C (chance error 0.875).

339 the bigger gains under the OVO strategy. Probably the OVO strategy can  
 340 filter some of the noisy features more efficiently than OVA, as it considers  
 341 significantly more lists of features ( $M = c(c - 1)/2$  vs.  $M = c$ ). After this  
 342 comparison we will only use OVO-SVM to evaluate our new methods.

### 343 3.4. Evaluation of the Methods on Artificial Datasets

344 Figure 4 shows the results for 4 versions of the Artificial-1 problem and 2  
 345 of the Artificial-3 problem. The remaining dataset from Artificial-1 is shown  
 346 on Panel (c) of Figure 3. Results for Artificial-2 are shown on Panel (d) of  
 347 the same figure. Additional datasets are included in the Additional Material  
 348 section. Overall, artificial problems show that the K-First method is the



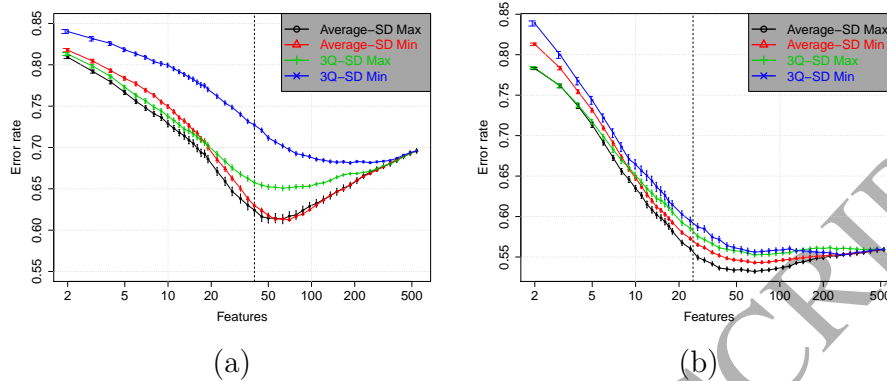


Figure 2: Comparison of tie breaking using maximum or minimum SD for Average-SD and 3Q-SD. Details are similar to Figure 1. Chance error of 0.875 for both panels. (a) Artificial-1-8C. (b) Artificial-2-8C.

349 most efficient one in finding subsets of features with low classification error,  
 350 followed closely by the Best-Ranking method. The Schulze method shows  
 351 good performance in several datasets. The other 3 methods show similar  
 352 results, though not as good as the first group.

353 On the Artificial-3 datasets (all relevant features are useful for all classes)  
 354 the differences among methods are clearly smaller than on the other 2 prob-  
 355 lems (with class-dependent features). Differences in performance increase  
 356 with the number of classes for the Artificial-1 dataset.

357 Comparing the two methods based on voting theory, the low performance  
 358 of Condorcet compared with Schulze is notorious. Taking a closer look at the  
 359 method, we noticed that Condorcet produces a lot of ties in the rankings,  
 360 which are broken using average positions. This produces a bias towards  
 361 features with good global average values instead of features highly relevant  
 362 for a few lists.

363 Another interesting analysis that can be made with artificial datasets is  
 364 the position occupied by the truly relevant features on the rankings produced  
 365 by the different methods, as we know in advance which features are noisy  
 366 and which ones are informative. A perfect method should rank all relevant  
 367 features first, with all noisy features following.

368 For each artificial problem, we analyzed the distribution of rankings given  
 369 by each selection method to the set of relevant features and to the set of  
 370 noisy features. We then computed some descriptive statistics of those two

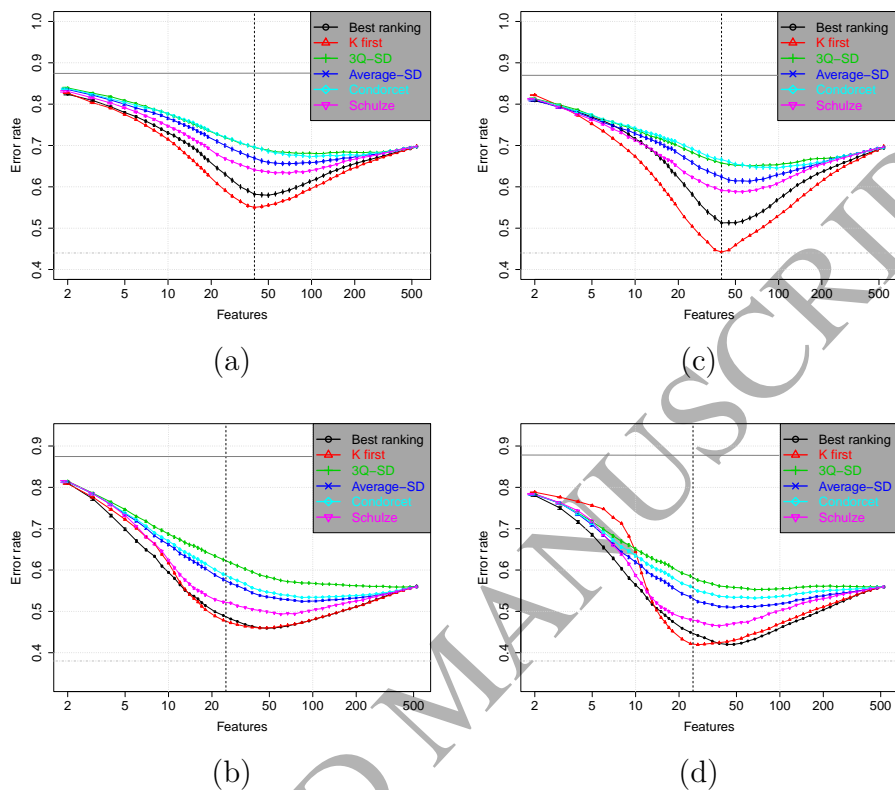


Figure 3: OVA-SVM vs. OVO-SVM on two artificial problems. Details are similar to Figure 1. (a) RFE-OVA-SVM on Artificial-1-8C. (b) RFE-OVA-SVM on Artificial-2-8C. (c) RFE-OVO-SVM on Artificial-1-8C (d) RFE-OVO-SVM on Artificial-2-8C.

371 distributions (Best, 1st. quartile, Mean, 3rd. quartile and Worst). In Table  
 372 4 we show these statistics on dataset Artificial-1-8C, which is representative  
 373 of the results obtained on the other versions of this problem. All six methods  
 374 rank relevant features at the first positions and noisy features at the last ones,  
 375 but there are important differences. Looking at the Mean and 3rd. Quartile  
 376 of the distributions, it is clear that K-First, Best Ranking and Schulze, in that  
 377 order, are the most accurate ones in ranking most of the features according  
 378 to their global relevance. These results confirm that the low error rates on  
 379 the figures discussed before are directly related to a better feature selection  
 380 by those methods.

381 In Table 5 we show the corresponding statistics for the Artificial-2-8C

Relevant	3Q-SD	Av-SD	Best Rank.	K-First	Condorcet	Schulze
Best	1	1	1	1	1	1
1st Q.	14	14	12	11	17	13
Mean	70	65	33	22	79	48
3rd Q.	98	86	41	31	110	61
Worst	436	476	357	495	474	410
Noisy	3Q-SD	Av-SD	Best Rank.	K-First	Condorcet	Schulze
Best	1	1	1	3	1	1
1st Q.	160	161	165	165	159	163
Mean	287	287	290	290	286	288
3rd Q.	415	415	415	415	415	415
Worst	540	540	540	540	540	540

Table 4: Statistics of the rankings given to relevant and noisy features by the diverse methods considered in this study for the Artificial-1-8C dataset. Values are rounded when needed.

dataset. This is the most interesting problem, as it contains subsets of relevant features with diverse levels of relevance. In the table we separated the relevant features into 3 subsets. As it can be observed in the table, K-First is the most effective strategy in separating the subsets of relevant features, and not only relevant from noisy features.

Finally, in Table 6 we show statistics for the Artificial-3-8C dataset — it is representative of all versions of this problem. As discussed before, all methods are almost equivalent on this dataset.

A last comment is in order about the Best Ranking method. As it can be seen in the tables, it can give high rankings to noisy features more easily than the K-First method, as it bases the final ranking on a single value for each feature.

#### 4. Evaluation on real-world datasets

We used 14 real-world datasets to evaluate our new methods. Details and origins of the datasets are collected on Table 7. We selected datasets from three different domains. The first 4 datasets collect mass spectrometry measurements of food products. All recorded peaks are present in the datasets. Some of the products under analysis can present class-specific features, reflecting particularities of some products, such as origin or manufacturing

1 to 5	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	1	1	1	1	1	1
1st Q.	2	2	2	2	2	2
Mean	22	5	9	3	15	4
3rd Q.	28	5	17	4	18	5
Worst	481	62	45	7	85	18
6 to 10	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	1	1	1	1	1	1
1st Q.	17	7	6	7	13	6
Mean	80	24	13	8	50	9
3rd Q.	158	30	19	9	58	11
Worst	523	241	44	14	453	70
11 to 25	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	1	2	1	8	1	4
1st Q.	35	19	12	14	29	15
Mean	161	87	22	18	101	65
3rd Q.	265	116	30	22	185	70
Worst	524	500	57	29	524	501
Noisy	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	1	3	3	10	1	8
1st Q.	142	147	151	151	142	148
Mean	269	273	275	275	274	274
3rd Q.	399	400	400	400	399	400
Worst	525	525	525	525	525	525

Table 5: Statistics of the rankings given to relevant and noisy features by the diverse methods considered in this study for the Artificial-2-8C dataset. The relevant features are divided into three subsets, and ordered according to their relevance by construction. Values are rounded when needed.

Relevant	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	1	1	1	1	1	1
1st Q.	11	11	11	11	11	11
Mean	27	29	41	26	28	29
3rd Q.	33	37	53	33	33	37
Worst	440	431	368	250	342	252
Noisy	3Q-SD	Av-SD	Best Rank.	K First	Condorcet	Schulze
Best	5	2	2	8	3	7
1st Q.	165	165	164	165	165	165
Mean	290	290	290	290	290	290
3rd Q.	415	415	415	415	415	415
Worst	540	540	540	540	540	540

Table 6: Statistics of the rankings given to relevant and noisy features by the diverse methods considered in this study for the Artificial-3-8C dataset. Values are rounded when needed.

401 method. The following 6 datasets come from the UCI repository. These  
 402 are more traditional datasets, with more samples than features and multiple  
 403 classes, involving typical pattern recognition problems. Finally, we selected  
 404 4 gene expression datasets from human tissues. These datasets were filtered  
 405 by curators to obtain circa 1000 genes with high signal-to-noise ratio in each  
 406 case.

407 In order to compare our results against previous methods we implemented  
 408 3 versions of MSVM-RFE, as described by Zhou & Tuck (2007). The first  
 409 method uses the multiclass SVM developed by Crammer & Singer (2001).  
 410 We will denote it "Zhou C & S". The second one uses the method of Weston  
 411 & Watkins (1999), in the following denoted by "Zhou W & W". Finally, we  
 412 implemented MSVM-RFE with the OVO decomposition of the traditional  
 413 binary SVM. We will refer to this method as "Zhou OVO". Notice that it  
 414 is equivalent to the Average Weights methodology, which is implemented by  
 415 default in most available Machine Learning packages, as previously explained.

416 On Figures 5 and 6 we show the results for eight datasets, while the  
 417 remaining cases are shown in the Additional Material section. In general,  
 418 differences in results for real world data are less notorious than for the Ar-  
 419 tificial problems. For UCI and Mass-Spectrometry datasets, K-First is in  
 420 general the method showing the best results in finding small subsets with  
 421 reduced classification error, followed by Best Ranking and 3rd Quartile. In

422 some problems, like Apples and Libra, differences are more notorious. On the  
423 gene expression datasets all methods show small differences, but in general  
424 the variants of MSVM-RFE exhibit better results than on the other domains.  
425 These datasets have been filtered by curators and as a consequence all fea-  
426 tures are informative. We believe that this improves the performance of  
427 averaging methods over methods that search for class-specific features.

428 Error curves show the complete behavior of the methods as a function  
429 of the number of features, but occasionally diverse methods are more effi-  
430 cient in selecting a high number or just a few features. To produce a more  
431 concrete comparison, we measured for two fixed numbers of selected features  
432 (10 and 20) the proportion of runs on which method A shows a smaller error  
433 than method B, for each of the three domains under evaluation. The full  
434 resulting matrices are shown in the Additional Material section. From these  
435 matrices we computed a ranking for each method, counting the number of  
436 other methods that it excels. We show the corresponding results in Table 8.  
437 They confirm the information extracted from the error curves: on the UCI  
438 and Mass-Spectrometry domains K-first shows the best results, but Best  
439 Ranking and 3rd Quartile also have high rankings. On the gene expression  
440 domains the best results come from one of the MSVM-RFE methods.

## 441 5. Computational burden

442 We evaluated the burden of the 6 new methods as a function of the  
443 number of features and samples using the Artificial-1 dataset. In panel (a)  
444 of Figure 7 we show how the running time scales with the number of features  
445 in the problems, using a log-scale for times. We include the 3 versions of the  
446 method by Zhou et al. as a comparison. It is clear from this figure that all but  
447 the Schulze method scale almost linearly with the number of features, being  
448 Condorcet and 3rd Quartile the slowest methods. Schulze is cubic in the  
449 number of features, as it involves a variant of the FloydWarshall algorithm  
450 to find shortest paths in a graph. On panel (b) of the same figure we show  
451 the dependence on the number of samples in the dataset. All new methods,  
452 including Schulze, scale almost linearly with the number of samples. The  
453 two variants of Zhou’s method using direct multiclass SVMs show power-  
454 law scaling with the number of samples.

## 455 6. Conclusions

456 In this work we discussed in depth the use of combinations of lists of  
457 features (instead of averaging individual importances) in SVM-RFE for fea-  
458 ture selection on multiclass problems. Using an appropriate mathematical  
459 framework we introduced 6 different methods to produce the final ranking of  
460 features starting from a set of ranked features list produced by each binary  
461 problem. We evaluated them in a series of artificial and real world datasets.

462 Our first conclusion is that the OVO strategy should be preferred over  
463 OVA for multiclass feature selection. Probably the higher number of binary  
464 problems in OVO helps in filtering out some noisy features that receive high  
465 rankings from just one or only a few binary problems, a similar beneficial  
466 effect to the use of ensembles in general.

467 Our second conclusion is that, overall, the K-First method is the most  
468 consistent one in selecting subsets of relevant features that lead to smaller  
469 classification errors. The idea is well-known in the document retrieval liter-  
470 ature, only considers the top k values of each list, and adapts efficiently to  
471 feature selection. We showed with several artificial and real-world datasets  
472 that this new method is superior to the typical weights averaging that is  
473 implemented by default in all current Machine Learning libraries.

474 Finally, two other methods also showed good results but present some  
475 drawbacks. The Best Ranking strategy is simple and efficient, but can lose  
476 performance on some problems, such as Artificial-2. Also, the use of a single  
477 value to characterize the behavior of a feature can give high rankings to noisy  
478 features by chance. The Schulze strategy, based on voting theory, shows a  
479 very good performance on some artificial datasets but does not compare well  
480 on real-world ones, and is by far the most complex and time-consuming  
481 strategy out of the six methods under evaluation.

482 Overall, the new methods were designed for problems with class-specific  
483 features, which is where they show their best performance. As they employ  
484 the OVO strategy, they are also resistant to noisy features. Filtered domains,  
485 with lots of low-relevance features and little noise like our gene expression  
486 datasets, seem to represent a more challenging domain for our new methods.

487 Work in progress includes a more extensive evaluation and the use of a  
488 penalty term to help discard correlated features.

489 **References**

- 490 Allwein, E. L., Schapire, R. E., & Singer, Y. (2000). Reducing multiclass  
491 to binary: A unifying approach for margin classifiers. *Journal of machine*  
492 *learning research*, *1*, 113–141.
- 493 Breiman, L. (2001). Random forests. *Machine learning*, *45*, 5–32.
- 494 Cappellin, L., Soukoulis, C., Aprea, E., Granitto, P., Dallabetta, N., Costa,  
495 F., Viola, R., Märk, T. D., Gasperi, F., & Biasioli, F. (2012). Ptr-  
496 tof-ms and data mining methods: a new tool for fruit metabolomics.  
497 *Metabolomics*, *8*, 761–770.
- 498 Crammer, K., & Singer, Y. (2001). On the algorithmic implementation of  
499 multiclass kernel-based vector machines. *Journal of machine learning re-*  
500 *search*, *2*, 265–292.
- 501 Del Pulgar, J. S., Soukoulis, C., Carrapiso, A., Cappellin, L., Granitto, P.,  
502 Aprea, E., Romano, A., Gasperi, F., & Biasioli, F. (2013). Effect of the  
503 pig rearing system on the final volatile profile of iberian dry-cured ham as  
504 detected by ptr-tof-ms. *Meat science*, *93*, 420–428.
- 505 Dittman, D. J., Khoshgoftaar, T. M., Wald, R., & Napolitano, A. (2013).  
506 Classification performance of rank aggregation techniques for ensemble  
507 gene selection. In *FLAIRS Conference*.
- 508 Duan, K.-B., Rajapakse, J. C., & Nguyen, M. N. (2007). One-versus-one and  
509 one-versus-all multiclass svm-rfe for gene selection in cancer classification.  
510 In *European Conference on Evolutionary Computation, Machine Learning*  
511 *and Data Mining in Bioinformatics* (pp. 47–56). Springer.
- 512 Fabris, A., Biasioli, F., Granitto, P. M., Aprea, E., Cappellin, L., Schuhfried,  
513 E., Soukoulis, C., Märk, T. D., Gasperi, F., & Endrizzi, I. (2010). Ptr-tof-  
514 ms and data-mining methods for rapid characterisation of agro-industrial  
515 samples: influence of milk storage conditions on the volatile compounds  
516 profile of trentingrana cheese. *Journal of mass spectrometry*, *45*, 1065–  
517 1074.
- 518 Forman, G. (2003). An extensive empirical study of feature selection metrics  
519 for text classification. *Journal of machine learning research*, *3*, 1289–1305.

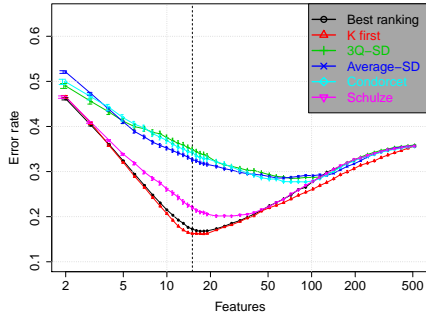


- 520 Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M.,  
521 Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A.  
522 et al. (1999). Molecular classification of cancer: class discovery and class  
523 prediction by gene expression monitoring. *science*, *286*, 531–537.
- 524 Granitto, P. M., Furlanello, C., Biasioli, F., & Gasperi, F. (2006). Recursive  
525 feature elimination with random forest for ptr-ms analysis of agroindustrial  
526 products. *Chemometrics and Intelligent Laboratory Systems*, *83*, 83–90.
- 527 Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature  
528 selection. *Journal of Machine Learning Research*, *3*, 1157–1182.
- 529 Guyon, I., Weston, J., Barnhill, S., & Vapnik, V. (2002). Gene selection for  
530 cancer classification using support vector machines. *Machine learning*, *46*,  
531 389–422.
- 532 Haury, A.-C., Gestraud, P., & Vert, J.-P. (2011). The influence of feature  
533 selection methods on accuracy, stability and interpretability of molecular  
534 signatures. *PloS one*, *6*, e28210.
- 535 Ho, T. K., Hull, J. J., & Srikari, S. N. (1994). Decision combination in  
536 multiple classifier systems. *IEEE transactions on pattern analysis and  
537 machine intelligence*, *16*, 66–75.
- 538 Hsu, C.-W., & Lin, C.-J. (2002). A comparison of methods for multiclass  
539 support vector machines. *IEEE transactions on Neural Networks*, *13*,  
540 415–425.
- 541 Hua, J., Tembe, W. D., & Dougherty, E. R. (2009). Performance of feature-  
542 selection methods in the classification of high-dimension data. *Pattern  
543 Recognition*, *42*, 409–424.
- 544 Jurman, G., Merler, S., Barla, A., Paoli, S., Galea, A., & Furlanello, C.  
545 (2008). Algebraic stability indicators for ranked lists in molecular profiling.  
546 *Bioinformatics*, *24*, 258–264.
- 547 Kanth, N. R., & Saraswathi, S. (2015). Efficient speech emotion recognition  
548 using binary support vector machines & multiclass svm. In *Computational  
549 Intelligence and Computing Research (ICIC)*, 2015 IEEE International  
550 Conference on (pp. 1–6). IEEE.

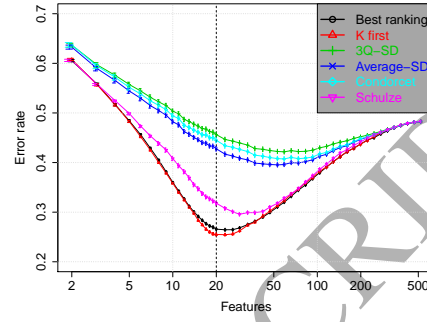
- 551 Kira, K., & Rendell, L. A. (1992). The feature selection problem: Traditional  
552 methods and a new algorithm. In *AAAI* (pp. 129–134). volume 2.
- 553 Kohavi, R., & John, G. H. (1997). Wrappers for feature subset selection.  
554 *Artificial intelligence*, *97*, 273–324.
- 555 Leek, J. T., Scharpf, R. B., Bravo, H. C., Simcha, D., Langmead, B., Johnson,  
556 W. E., Geman, D., Baggerly, K., & Irizarry, R. A. (2010). Tackling the  
557 widespread and critical impact of batch effects in high-throughput data.  
558 *Nature Reviews Genetics*, *11*, 733–739.
- 559 Lichman, M. (2013). UCI machine learning repository.
- 560 Liu, H., Dougherty, E. R., Dy, J. G., Torkkola, K., Tuv, E., Peng, H., Ding,  
561 C., Long, F., Berens, M., Parsons, L. et al. (2005). Evolving feature selec-  
562 tion. *IEEE Intelligent systems*, *20*, 64–76.
- 563 Mohsenzadeh, Y., Sheikhzadeh, H., & Nazari, S. (2016). Incremental rel-  
564 evance sample-feature machine: A fast marginal likelihood maximization  
565 approach for joint feature selection and classification. *Pattern Recognition*,  
566 *60*, 835–848.
- 567 Mohsenzadeh, Y., Sheikhzadeh, H., Reza, A. M., Bathaee, N., & Kalayeh,  
568 M. M. (2013). The relevance sample-feature machine: A sparse bayesian  
569 learning approach to joint feature-sample selection. *IEEE transactions on*  
570 *cybernetics*, *43*, 2241–2254.
- 571 Monti, S., Tamayo, P., Mesirov, J. P., & Golub, T. R. (2003). Consensus clus-  
572 tering: A resampling-based method for class discovery and visualization of  
573 gene expression microarray data. *Machine Learning*, *52*, 91–118.
- 574 Neumayer, R., Mayer, R., & Nørnvåg, K. (2011). Combination of feature  
575 selection methods for text categorisation. In *European Conference on In-*  
576 *formation Retrieval* (pp. 763–766). Springer.
- 577 Nguyen, M. H., & De la Torre, F. (2010). Optimal feature selection for  
578 support vector machines. *Pattern recognition*, *43*, 584–591.
- 579 Nuray, R., & Can, F. (2006). Automatic ranking of information retrieval  
580 systems using data fusion. *Information processing & management*, *42*,  
581 595–614.

- 582 Raies, A. B., & Bajic, V. B. (2016). In silico toxicology: computational  
583 methods for the prediction of chemical toxicity. *Wiley Interdisciplinary*  
584 *Reviews: Computational Molecular Science*, 6, 147–172.
- 585 Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C.-H., An-  
586 gelo, M., Ladd, C., Reich, M., Latulippe, E., Mesirov, J. P. et al. (2001).  
587 Multiclass cancer diagnosis using tumor gene expression signatures. *Pro-*  
588 *ceedings of the National Academy of Sciences*, 98, 15149–15154.
- 589 Saari, D. (2001). *Chaotic elections!: A mathematician looks at voting*. Amer-  
590 ican Mathematical Soc.
- 591 Schulze, M. (2011). A new monotonic, clone-independent, reversal symmet-  
592 ric, and condorcet-consistent single-winner election method. *Social Choice*  
593 *and Welfare*, 36, 267–303.
- 594 Statnikov, A., Aliferis, C. F., Tsamardinos, I., Hardin, D., & Levy, S. (2005).  
595 A comprehensive evaluation of multicategory classification methods for  
596 microarray gene expression cancer diagnosis. *Bioinformatics*, 21, 631–643.
- 597 Uysal, A. K. (2016). An improved global feature selection scheme for text  
598 classification. *Expert systems with Applications*, 43, 82–92.
- 599 Vapnik, V. (2013). *The nature of statistical learning theory*. Springer science  
600 & business media.
- 601 Weston, J., Mukherjee, S., Chapelle, O., Pontil, M., Poggio, T., & Vapnik, V.  
602 (2000). Feature selection for svms. In *Proceedings of the 13th International*  
603 *Conference on Neural Information Processing Systems* (pp. 647–653). MIT  
604 Press.
- 605 Weston, J., & Watkins, C. (1999). Support vector machines for multi-class  
606 pattern recognition. In *ESANN* (pp. 219–224). volume 99.
- 607 You, W., Yang, Z., & Ji, G. (2014). Feature selection for high-dimensional  
608 multi-category data using pls-based local recursive feature elimination. *Ex-*  
609 *pert Systems with Applications*, 41, 1463–1475.
- 610 Young, H. P. (1988). Condorcet’s theory of voting. *The American Political*  
611 *Science Review*, (pp. 1231–1244).

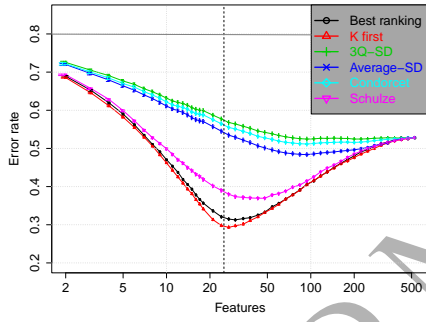
- 612 Zhou, N., & Wang, L. (2016). Processing bio-medical data with class-  
613 dependent feature selection. In *Advances in Neural Networks* (pp. 303–  
614 310). Springer.
- 615 Zhou, W., & Dickerson, J. A. (2014). A novel class dependent feature se-  
616 lection method for cancer biomarker discovery. *Computers in biology and*  
617 *medicine*, *47*, 66–75.
- 618 Zhou, X., & Tuck, D. P. (2007). Msvm-rfe: extensions of svm-rfe for mul-  
619 ticlass gene selection on dna microarray data. *Bioinformatics*, *23*, 1106–  
620 1114.



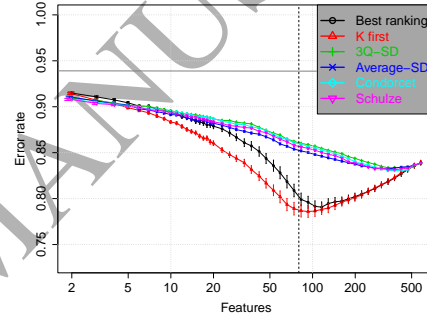
(a)



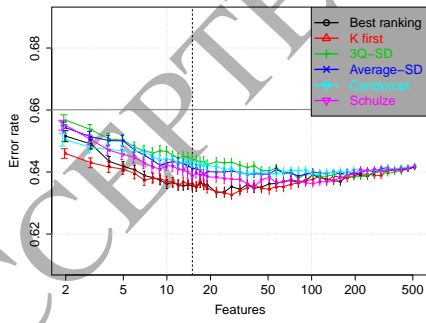
(b)



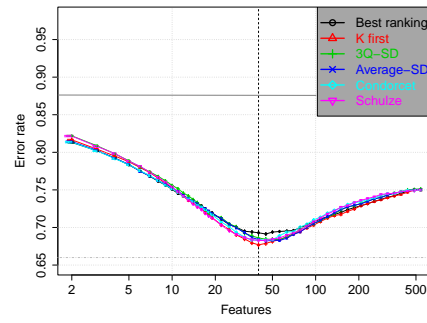
(c)



(d)



(e)



(f)

Figure 4: Error curves for the six methods on some artificial datasets. Details are similar to Figure 1. (a) Artificial-1-3C (chance error 0.666) (b) Artificial-1-4C (chance error 0.75) (c) Artificial-1-5C (d) Artificial-1-16C (e) Artificial-3-3C (f) Artificial-3-8C

Dataset	F	S	C	Description
Apple	714	150	15	Mass spectrometry measurements over 15 varieties of Apple clones (Cappellin et al., 2012)
Cheese	117	72	8	Mass spectrometry measurements over 8 varieties of Italian cheese (Fabris et al., 2010)
Ham	1338	382	11	Mass spectrometry measurements over 11 varieties of Iberian Hams (Del Pulgar et al., 2013)
Strawberry	232	233	9	Mass spectrometry measurements over 9 varieties of strawberries (Granitto et al., 2006)
Multi-F	649	2000	10	Features of handwritten numerals extracted from a collection of Dutch utility maps (Lichman, 2013)
Libras	90	360	15	Diverse hand movements from the Brazilian hands language (Lichman, 2013)
Robot1	90	88	4	Robot Execution Failures Data Set, from UCL. Failures in approach to grasp position (Lichman, 2013)
Robot3	90	47	4	Same as Robot1. Position of part after a transfer failure
Robot4	90	117	3	Same as Robot1. Failures in approach to ungrasp position
Robot5	90	164	5	Same as Robot1. Failures in motion with part
Leukemia	985	248	6	Gene expression of Bone marrow samples with 6 subtypes of Leukemia (Monti et al., 2003)
Lung	1000	197	4	Gene expression of lung tissues with 4 cancer types (Monti et al., 2003)
CNS	989	42	5	Gene expression of 5 tumor types of the central nervous system (Monti et al., 2003)
Novartis	1000	103	4	Gene expression of tissue samples from 4 distinct cancer types (Monti et al., 2003)

Table 7: Details on the 14 real-world datasets used in this work. Columns show the number of features (F), samples (S) and classes (C).

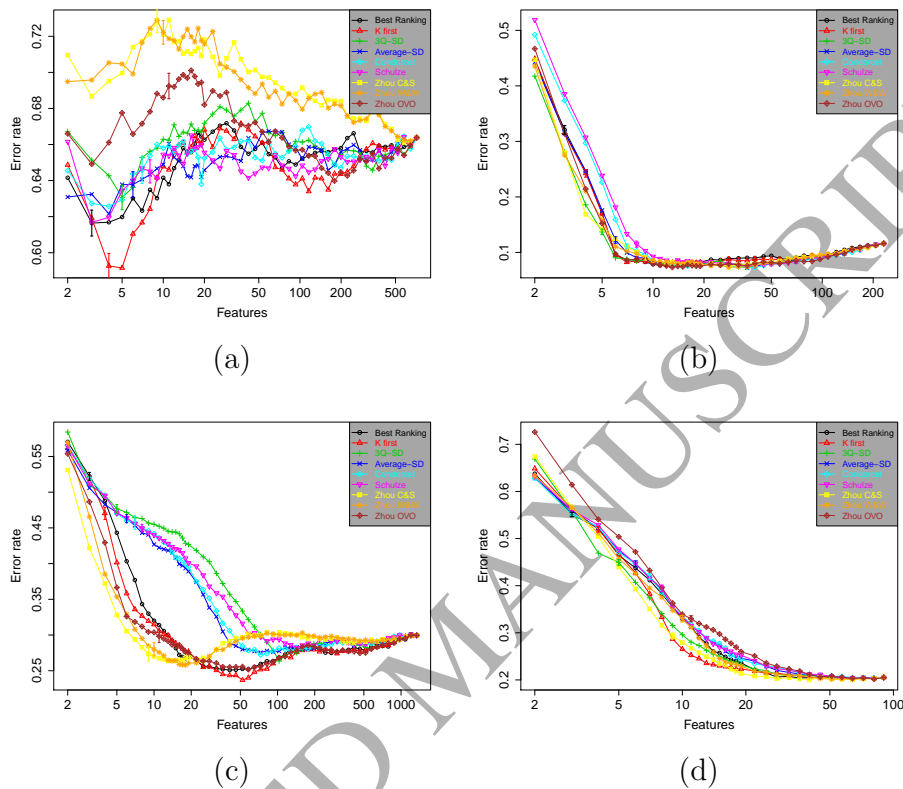


Figure 5: Results on some real world datasets. Details are similar to Figure 1. (a) Apples. (b) Strawberry. (c) Ham. (d) Libras.

Group	3Q-SD	Av-SD	Best	K-First	Cond.	Schulze	Z. C&S	Z. W&W	Z. OVO
UCI-10 Feat	7	5	5	<b>8</b>	3	2	0	2	5
MS -10 Feat	0	3	7	<b>8</b>	2	1	5	4	6
GE -10 Feat	1	3	6	4	2	0	7	5	<b>8</b>
UCI-20 Feat	5	2	6	<b>8</b>	2	0	3	7	3
MS -20 Feat	0	2	7	<b>8</b>	3	1	4	5	6
GE -20 Feat	1	6	5	4	2	0	<b>8</b>	3	7

Table 8: Rankings of methods (higher is better) counting the number of times that one method outperforms another on each domain and number of selected features.

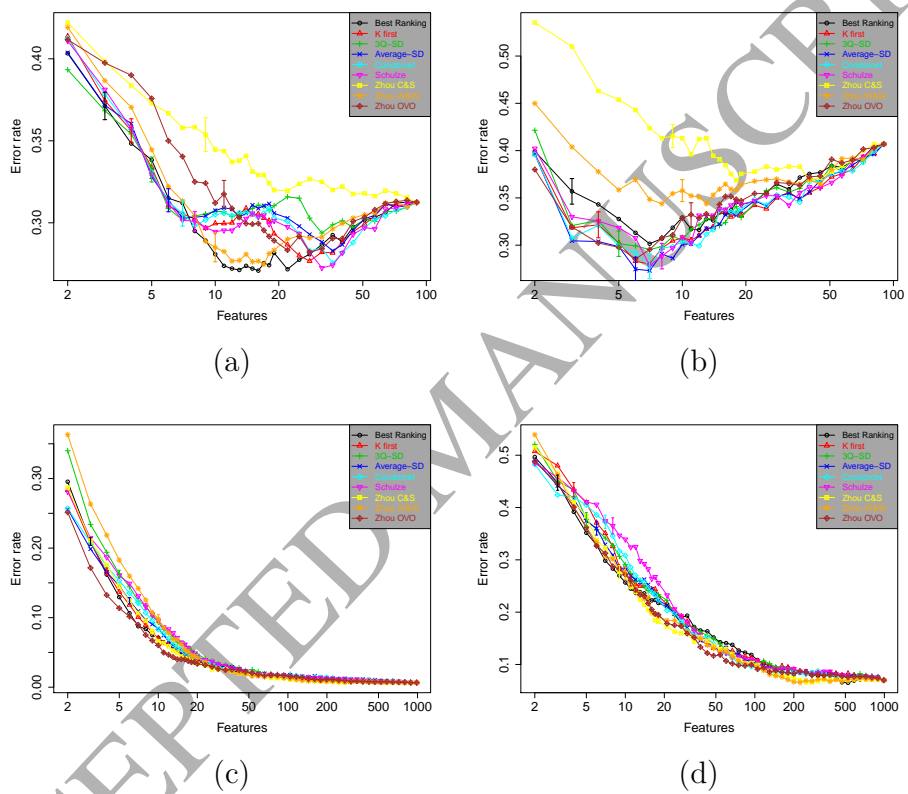


Figure 6: Results on some real world datasets. Details are similar to Figure 1. (a) Robot1. (b) Robot3. (c) Leukemia. (d) CNS.



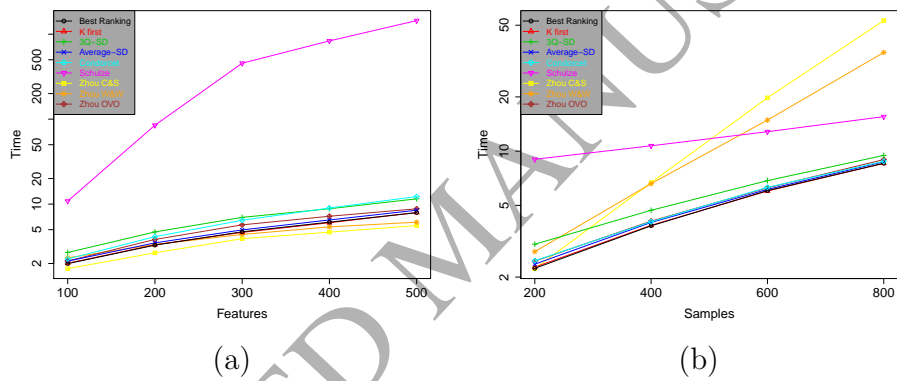


Figure 7: Comparison of running times for all methods evaluated in this work as a function of (a) the number of features and (b) the number of samples. Times (in seconds) are in log-scale.