

# Error curves for evaluating the quality of feature rankings

Ivica Slavkov<sup>1</sup>, Matej Petković<sup>Corresp., 1, 2</sup>, Pierre Geurts<sup>3</sup>, Dragi Kocev<sup>1, 2</sup>, Sašo Džeroski<sup>1, 2</sup>

<sup>1</sup> Jozef Stefan Institute, Ljubljana, Slovenia

<sup>2</sup> Jozef Stefan International Postgraduate School, Ljubljana, Slovenia

<sup>3</sup> Université de Liège, Liège, Belgium

Corresponding Author: Matej Petković

Email address: matej.petkovic@ijs.si

In this paper, we propose a method for evaluating feature ranking algorithms. A feature ranking algorithm estimates the importance of descriptive features when predicting the target variable, and the proposed method evaluates the correctness of these importance values by computing the error measures of two chains of predictive models. The models in the first chain are built on nested sets of top-ranked features, while the models in the other chain are built on nested sets of bottom ranked features. We investigate which predictive models are appropriate for building these chains, showing empirically that the proposed method gives meaningful results and can detect differences in feature ranking quality. This is first demonstrated on synthetic data, and then on several real-world classification benchmark problems.

# Error curves for evaluating the quality of feature rankings

Ivica Slavkov<sup>1</sup>, Matej Petković<sup>1,2</sup>, Pierre Geurts<sup>3</sup>, Dragi Kocev<sup>1,2</sup>, and Sašo Džeroski<sup>1,2</sup>

<sup>1</sup>Jožef Stefan Institute, Ljubljana, Slovenia

<sup>2</sup>Jožef Stefan International Postgraduate School, Ljubljana, Slovenia

<sup>3</sup>University of Liège, Liège, Belgium

Corresponding author:

Matej Petković<sup>1,2</sup>

Email address: matej.petkovic@ijs.si

## ABSTRACT

In this paper, we propose a method for evaluating feature ranking algorithms. A feature ranking algorithm estimates the importance of descriptive features when predicting the target variable, and the proposed method evaluates the correctness of these importance values by computing the error measures of two chains of predictive models. The models in the first chain are built on nested sets of top-ranked features, while the models in the other chain are built on nested sets of bottom ranked features. We investigate which predictive models are appropriate for building these chains, showing empirically that the proposed method gives meaningful results and can detect differences in feature ranking quality. This is first demonstrated on synthetic data, and then on several real-world classification benchmark problems.

## INTRODUCTION

In the era of data abundance, we face high-dimensional problems increasingly often. Sometimes, prior to applying predictive modeling (e.g., classification) algorithms to such problems, dimensionality reduction may be necessary for a number of reasons, including computational reasons. By keeping only a limited number of descriptors (features), a classifier can also achieve better predictive performance, since typically, a portion of the features strongly influence the target variable, and the others can be understood as (mostly) noise. This dimensionality reduction corresponds to the task of feature selection (Guyon et al., 2002). A task related to it is feature ranking. This is a generalization of feature selection where, in addition to simply telling apart relevant features from irrelevant ones (Nilsson et al., 2007), one also assesses how relevant are they for predicting the target variable.

In machine learning, feature ranking is typically seen either as a preprocessing or as a postprocessing step. In the former case, one actually tackles the feature selection problem by first computing the feature relevance values, and then keeping only the features whose relevance is above some user defined threshold. In the second case, feature ranking is obtained after building a predictive model in order to explain it, e.g., (Arceo-Vilas et al., 2020). For black box models, such as neural networks, this may be the only way to understand their predictions.

In some application domains, such as biology or medicine, feature ranking may be the main point of interest. If we are given data about the expression of numerous genes, for a group of patients, and the patients' clinical state (diseased/healthy), one can find good candidate genes that influence the health status of the patients, which gives us a deeper understanding of the disease.

Due to the prominence of the feature ranking task, there exist many feature ranking methods. Simpler methods assess the relevance of each feature independently ignoring the other features ( $\chi^2$  statistics, mutual information of the feature and the target variable) and their possible interactions. A prominent example that shows the myopic nature of such approaches is the case when the target variable  $y$  is defined as  $y = \text{XOR}(x_1, x_2)$  where  $x_1$  and  $x_2$  are two binary features. Ignoring  $x_1$  when computing the relevance of  $x_2$  (and vice-versa) would result in assessing  $x_1$  as completely irrelevant, i.e., as random noise. More

sophisticated methods assess relevance of each feature in the context of the others. They are typically based on some predictive model, e.g., Random Forest feature ranking (Breiman, 2001), or optimisation problem (Nardone et al., 2019), but not necessarily, cf. e.g., ReliefF (Robnik-Šikonja and Kononenko, 2003) and the work of Li and Gu (2015).

However, there is no unified definition of feature importance, and actually, every feature ranking algorithm comes with its own (implicit) definition. Therefore, different methods typically introduce different feature importance scores: Deciding which of them is the best is a very relevant, but also very challenging task that we would like to address in this paper. More precisely, we continue and extend our previous work (Slavkov et al., 2018), where we proposed and evaluated a quantitative score for the assessment of feature rankings. Here, we propose a new feature ranking evaluation method that can evaluate feature rankings in a relative sense (deciding which of the feature rankings is better), or in an absolute sense (assessing how good is a feature ranking). The method is based on constructing two chains of predictive models that are built from the top-ranked and bottom-ranked features. The predictive performances of the models in the chain are then shown on graphs of so called forward feature addition (FFA) and reverse feature addition (RFA) curves, which reveal how the relevant features are distributed in the ranking(s). An important property of the proposed method is that it does not need any prior ground truth knowledge of the data.

We investigate the performance of the proposed evaluation approach under a range of scenarios. To begin with, we prove the potential of the FFA and RFA curves by using them in setting which employs synthetic data. Next, we investigate the use of different types of predictive models for constructing the curves, thus considerably extending the preliminary experiments by Slavkov et al. (2018). Furthermore, we apply the proposed evaluation approach to a large collection of benchmark datasets. Compared to Slavkov et al. (2018), we have included 11 new high-dimensional datasets. The results of the evaluation, in a nutshell, show that the FFA and RFA curves are able to discern the best ranking among multiple proposed feature rankings.

The remainder of this paper is organized as follows. Sec. *Related Work* outlines related work, Sec. *A Method for Evaluating Feature Rankings* describes in detail the proposed method for constructing error curves. Next, Sec. *Empirical Evaluation of FFA/RFA Curve Properties* discusses the properties of the error curves when applied to synthetic data. We then give the results of the experimental evaluation on benchmark datasets in Sec. *Feature Ranking Comparison on Real-Worlds Datasets*. Sec. *Conclusions* concludes with a summary of our contributions and an outline of possible directions for further work. In the appendices, we give additional information about generating synthetic data (Appendix A1), measuring distance between rankings (Appendix A2), and comparative evaluation of feature ranking methods (Appendix A3). In Appendix A4, detailed experimental results are given.

## RELATED WORK

The evaluation of feature rankings is a complex and unsolved problem. Typically, feature rankings are evaluated on artificially generated problems, while evaluation on real world problems remains an issue approached indirectly. To begin with, when the ground truth ranking is known, one can transform the problem of feature ranking evaluation into an evaluation of classification predictive model (Jong et al., 2004) as follows. First, a ranking is computed. Then, for every threshold, the numbers of relevant features (true positives) and irrelevant features (false positives) with the feature relevance above the threshold are computed. From these values, a ROC curve can be created and the area underneath it computed.

Another possible approach is to compute separability (Robnik-Šikonja and Kononenko, 2003), i.e., the minimal difference between the feature importance of a relevant feature and the feature importance of an irrelevant feature. If this difference is positive, then the relevant features are separated from the irrelevant ones, otherwise they are mixed.

However, both approaches are more applicable to feature selection problems and are too coarse for feature rankings problem, since they only differentiate between relevant and irrelevant features. Spearman's rank correlation coefficient between the computed and the ground truth ranking might be more appropriate.

The main shortcoming of the upper approaches is that they demand the ground truth ranking. In real world scenarios, this is not known, which makes the upper approaches useless. Nevertheless, using synthetic data and the controlled environment offers a good starting point for showing the usefulness of a feature ranking evaluation method, as we shall also see later.

An approach that overcomes the issue of unknown ground truth ranking bases on selecting  $k$  top-ranked features and building a predictive model that uses only these features to predict target variable. The ranking whose top-ranked features result in the model with the highest predictive performance, is proclaimed the best. Since it is now always clear which value of  $k$  should be chosen, this can be done for multiple values of  $k$  (Guyon et al., 2002; Furlanello et al., 2003; Paoli et al., 2005; Verikas et al., 2011).

In addition to correctness, rankings stability is sometimes also part of the evaluation. The stability of a ranking algorithm can be measured by comparing the feature rankings obtained, for example, from the different bootstrap replicates of a dataset or from the folds in cross-validation (Guzmán-Martínez and Alaiz-Rodríguez, 2011; Kalousis et al., 2007; Jurman et al., 2008). In (Saeys et al., 2008) both stability and predictive performance are combined into a single feature ranking quality index.

Also, notions similar to FFA curves (without any particular name, though) as the feature ranking evaluation method can be found in the literature ((Liu et al., 2003), (Duch et al., 2004), (Biesiada et al., 2005) and (Liang et al., 2008)). However, to the best of our knowledge, there is no discussion and detailed investigation why FFA curves are an appropriate method for comparing feature rankings, nor which learning methods should (or should not) be used for constructing them.

## A METHOD FOR EVALUATING FEATURE RANKINGS

First of all, every feature ranking method should be able to tell apart relevant features from irrelevant ones (Nilsson et al., 2007). In addition to that, the method should order the features with respect to the target variable, awarding the most relevant ones the top ranks.

If ground truth ranking exists, the method should return this ranking in the optimal case. The worst case is more complicated and has two possible answers. One is the inverse of the ground truth ranking. However, since the ground truth ranking is typically not known in real-world scenarios, a more useful definition of the worst ranking is random ranking. This ranking also contains as little information about the distribution of the relevant features in the ranking as possible. Moreover, this distribution can be always assessed and is the cornerstone of our ranking evaluation method.

### The Evaluation Method

First, we define the notation used in the rest of the paper:  $\mathcal{D}$  denotes a dataset whose columns are input features  $F_i$  that form a set  $\mathcal{F}$ , and the target feature  $F_t$ . A feature ranking method takes the dataset as an input, and returns a list  $\mathbf{R} = (F_{(1)}, \dots, F_{(n)})$  as the output, where  $F_{(i)}$  is the feature with the rank  $i$ .

We evaluate a ranking  $\mathbf{R}$  by evaluating different subsets  $\mathcal{S}$  of features  $\mathcal{F}$ . This is done by building a predictive model  $\mathcal{M}(\mathcal{S}, F_t)$  and assessing its predictive power. The evaluation of the predictive model provides a cumulative assessment of the information contained in the feature set  $\mathcal{S}$  and it can be quantified with an error measure  $\text{err}(\mathcal{M}(\mathcal{S}, F_t))$ . The question is how to generate the feature subsets from the feature ranking, so that the error estimates can provide insight into the correctness of the feature ranking and constitute an evaluation thereof.

The construction of the feature subsets should be guided by the feature ranking  $\mathbf{R}$ . Starting from the top ranked feature  $F_{(1)}$  and going towards the bottom ranked feature  $F_{(n)}$ , the feature relevance should decrease. Following this basic intuition, we propose two methods for constructing feature subsets from the feature ranking: *forward feature addition* (FFA) and *reverse feature addition* (RFA).

FFA constructs the feature subsets  $\mathcal{S}^i$  by considering the  $i$  highest ranked features, starting with  $\mathcal{S}^1 = \{F_{(1)}\}$ . The next set  $\mathcal{S}^{i+1}$  is constructed by adding the next lower-ranked feature, namely  $\mathcal{S}^{i+1} = \mathcal{S}^i \cup \{F_{(i+1)}\}$ . The process continues until  $i = n$  and  $\mathcal{S}^n$  contains all of the features from  $\mathbf{R}$ .

RFA produces feature sets  $\mathcal{S}_i$  constructed in an opposite manner to FFA. We start with  $\mathcal{S}_1 = \{\mathcal{F}_{(n)}\}$  that contains only the lowest ranked feature. The next feature set  $\mathcal{S}_{i+1}$  is constructed by adding the lowest-ranked feature which is not in  $\mathcal{S}_i$ , namely  $\mathcal{S}_{i+1} = \mathcal{S}_i \cup \{F_{(n-i)}\}$ . In the same way as for FFA, the process of RFA continues until we include all of the features, i.e.,  $\mathcal{S}_n = \mathcal{F}$ .

Note that FFA can be viewed as backward feature elimination. Starting from  $\mathcal{S}^n = \mathcal{F}$ , at each step we remove the least relevant feature from  $\mathcal{S}^i$  to obtain  $\mathcal{S}^{i-1}$ . Similarly, RFA can be viewed as forward feature elimination. Finally, it holds that  $\mathcal{F} = \mathcal{S}^{n-i} \cup \mathcal{S}_i$  for all  $i$ .

For each  $i$  and each constructed feature subset  $\mathcal{S}^i$  or  $\mathcal{S}_i$ , we build predictive models  $\mathcal{M}(\mathcal{S}^i, F_t)$  and  $\mathcal{M}(\mathcal{S}_i, F_t)$ . We then estimate their respective prediction errors,  $\text{err}^i$  and  $\text{err}_i$ . This results in two error curves. We name them FFA and RFA curves, each constructed by the corresponding FFA/RFA feature subset construction method. The value for each point of the FFA curve is defined as  $\text{FFA}(i) = \text{err}^i$ , while

---

**Algorithm 1** Generation of the FFA and RFA curves.

---

**Input:** Feature Ranking  $\mathbf{R} = (F_{(1)}, \dots, F_{(n)})$ , Target Feature  $F_t$ , type of curve (FFA or RFA)

$\mathcal{S} \leftarrow \emptyset$

$E \leftarrow \text{list of length } n$

**for**  $i = 1, 2, \dots, n$  **do**

**if** curve type is FFA **then**

$\mathcal{S} \leftarrow \mathcal{S} \cup \{F_{(i)}\}$

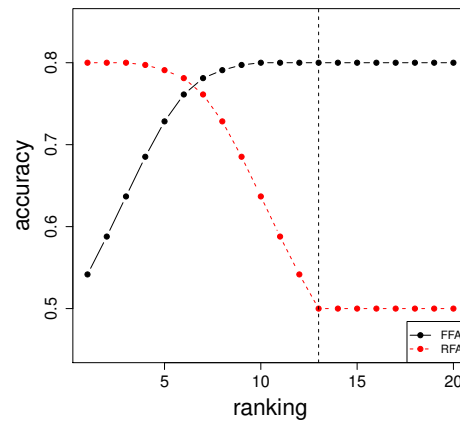
**else**

$\mathcal{S} \leftarrow \mathcal{S} \cup \{F_{(n-i+1)}\}$

$E[i] \leftarrow \text{err}(\mathcal{M}(\mathcal{S}, F_t))$

**return**  $E$

---



**Figure 1.** Sample FFA and RFA curves

153 for the RFA curve as  $\text{RFA}(i) = \text{err}_i$ . The process of FFA/RFA curve construction is summarised in  
 154 Algorithm 1.

155 The computational complexity of the proposed algorithm for constructing a single (FFA or RFA) curve  
 156 is  $\mathcal{O}(n(M + T))$ , where  $n$  is the number of features,  $M = M(n)$  is the cost of constructing the predictive  
 157 model and  $T = T(n)$  is the cost of its evaluation. It should be noted that  $M$  and  $T$  are dependent on the  
 158 specific learning method used for inducing the model and on the procedure used for evaluating it.

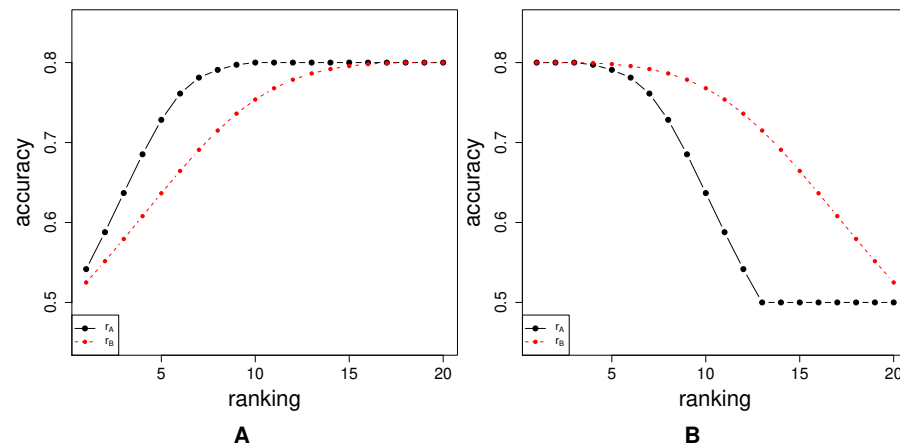
159 Typically, the points  $\text{FFA}(i)$ ,  $\text{FFA}(i+1)$ ,  $\cdot$ , do not differ considerably, for  $i$  large enough, since it  
 160 expected that only a small proportion of the features is relevant when the data is high-dimensional. This  
 161 means that we can make the algorithm more efficient if we construct the set  $\mathcal{S}^{i+\delta(i)}$  from the set  $\mathcal{S}^i$  by  
 162 including  $\delta(i)$  features into it. Analogously, we speed up the construction of the RFA curves.

### 163 Interpretation of Error Curves

164 The visualization and interpretation of the FFA and RFA curves can be explained by considering the  
 165 examples of FFA and RFA curves given in Fig. 1. The y-axis for both curves is the same and depicts the  
 166 error estimate of a feature subset. Point  $i$  at the x-axis corresponds to the moment when the feature  $F_{(i)}$   
 167 is first included in the predictive model:  $\mathcal{S}^i$  for the FFA curve and  $\mathcal{S}_{n-i+1}$  for the RFA curve. Thus, the  
 168 FFA curve in Fig. 1 is constructed from left-to-right as the top-ranked features are at the beginning of the  
 169 ranking. In contrast, the RFA curve is constructed from right-to-left starting with the end of the ranking  
 170 and going towards its beginning.

171 Let us first focus on the FFA curve. We can observe that as the number  $k$  of features increases, the  
 172 accuracy of the predictive models also increases. This can be interpreted as follows: By adding more  
 173 and more of the top- $k$  ranked features, the number of relevant features in the constructed feature subsets  
 174 increases, which is reflected in the improvement of the accuracy (error) measure.

175 Next, for the RFA curve in Fig. 1, if we inspect it from right-to-left, we can notice that it is quite  
 176 different from the FFA curve at the beginning. Namely, the accuracy of the models constructed with  
 177 the bottom ranked features is minimal, which means the ranking is correct in the sense that it puts only



**Figure 2.** Comparison of FFA curves (A) and RFA curves (B) of two ranking methods  $r_A$  and  $r_B$ .

irrelevant features at the bottom of the ranking. As the number of bottom- $k$  features increases, some relevant features are included and the accuracy of the models increases.

We now consider the complete Fig. 1. The FFA and RFA curve essentially provide an estimate of how the relevant features are spread throughout the feature ranking. Namely, the FFA curve provides us with an estimate of where the relevant features appear at the top of the ranking, while the RFA curve provides an estimate of where relevant features appear at the bottom of the ranking. In the specific case depicted in Fig. 1, the relevant features are located between the 1st and the 13th ranked feature.

Besides providing an estimate of the spread of the relevant features across the feature ranking, the real utility of the FFA/RFA curves becomes apparent if we consider them in a relative, or more precisely, a comparative context. Let us consider two arbitrary feature ranking methods  $r_A$  and  $r_B$ , which produce feature rankings  $\mathbf{R}_A$  and  $\mathbf{R}_B$ , respectively. For these two rankings, we present the corresponding FFA/RFA curves in Fig. 2.

We first inspect the FFA curves visually. We find that the values of the FFA curve of the ranking method  $r_A$  are (most of the time) above the FFA curve of the other ranking method  $r_B$ . This can be interpreted in the following way: for an arbitrary  $k$ , when considering the top- $k$  features of the feature rankings  $\mathbf{R}_A$  and  $\mathbf{R}_B$ , more relevant features are included in the top- $k$  features of ranking  $\mathbf{R}_A$  than the top- $k$  features of ranking  $\mathbf{R}_B$ . This implies that ranking algorithm  $r_A$  produces a better ranking as compared to the ranking algorithm  $r_B$ .

A similar discussion applies to the RFA curve. When one considers the bottom- $k$  features of a given feature ranking, most of the time, feature ranking  $\mathbf{R}_A$  includes less relevant features than feature ranking  $\mathbf{R}_B$ , i.e., the predictive models constructed are less accurate. Here, because the opposite logic of the FFA curve applies, one can also conclude that the feature ranking algorithm  $r_A$  produces a better feature ranking than the feature ranking algorithm  $r_B$ .

## Expected FFA and RFA Curves

When one wants to assess the quality of a single feature ranking in a real-world application, its forward (reverse) feature addition curves can be only compared to the curves that belong to the ranking, generated uniformly at random, since the ground-truth ranking is not known. As discussed before, the random ranking  $\mathbf{R}_{\text{RND}}$  is the worst-case ranking, since it contains no information about the distribution of the relevant features. As such, it can also serve as a baseline.

The expected values of the points that define FFA curve of the ranking  $\mathbf{R}_{\text{RND}}$  coincide with the expected values of the RFA curve of this ranking, since the corresponding values only depend on the data itself and the number of features  $i$  at a given point of the curves. Thus, *expected curves* can be the common name for both types of the curves that belong to  $\mathbf{R}_{\text{RND}}$ . Computing the exact average error estimations  $\mathbb{E}_{\mathcal{S}}[\text{err}_i^f] = \mathbb{E}_{\mathcal{S}}[\text{err}_i^r]$ , where  $\mathcal{S} \subseteq \mathcal{F}$ ,  $|\mathcal{S}| = i$ , may be unfeasible if the number of features  $n$  is large (e.g., for  $i = n/2$ ,  $\mathcal{O}((2n)!/(n!)^2)$  models have to be evaluated), but one can overcome this by sampling the sets  $\mathcal{S}$ .

## Stability of feature ranking

An important aspect of feature ranking algorithms is their stability (Nogueira et al., 2017) or, more specifically, the stability of the ranked feature lists that they produce. Once we have the set  $\mathcal{R}$  of  $m$  rankings  $\mathbf{R}_t$  that were induced from the different samples of the dataset  $\mathcal{D}$ , the stability index  $S(\mathcal{R})$  can be calculated as

$$St(\mathcal{R}) = \frac{1}{\binom{m}{2}} \sum_{t=1}^{m-1} \sum_{s=t+1}^m S_M(\mathbf{R}_t, \mathbf{R}_s),$$

i.e., the stability index is the average of pairwise similarities  $S_M$  for each pair of rankings. In general, the function  $S_M$  can be any (dis)similarity measure, e.g., the Spearman rank correlation coefficient (Saeys et al., 2008; Khoshgoftaar et al., 2013), the Jaccard distance (Saeys et al., 2008; Kalousis et al., 2007), an adaptation of the Tanimoto distance (Kalousis et al., 2007), Fuzzy (Goodman and Kruskal's) gamma coefficient (Boucheham and Batouche, 2014; Henzgen and Hüllermeier, 2015), etc.

To assess the stability of feature ranking in our experimental work, we set  $S_M = \text{Ca}$ , where Ca is the Canberra distance (Lance and Williams, 1966, 1967; Jurman et al., 2008). This is a weighted distance metric that puts bigger emphasis on the stability of the top ranked features. If we have two feature rankings  $\mathbf{R}_A$  and  $\mathbf{R}_B$  of  $n$  features, then the Canberra distance is defined as

$$\text{Ca}(\mathbf{R}_A, \mathbf{R}_B) = \sum_{j=1}^n \frac{|rank_A(F_j) - rank_B(F_j)|}{rank_A(F_j) + rank_B(F_j)}. \quad (1)$$

However, we do not only estimate the stability of the ranking as a whole. Rather, we also estimate the stability of the partial rankings based on the features from  $\mathcal{S}^i$ . In order for the distance to be applicable to such partial rankings with  $i < n$  features, the following adaptation is proposed: instead of using the ranks  $rank_{A,B}(F)$ , we use  $rank_{A,B}^i(F) = \min\{rank_{A,B}(F), i+1\}$ , i.e., all features with rank higher than  $i$  are treated as if they had rank  $i+1$ :

$$\text{Ca}(\mathbf{R}_A, \mathbf{R}_B) = \sum_{j=1}^n \frac{|rank_A^i(F_j) - rank_B^i(F_j)|}{rank_A^i(F_j) + rank_B^i(F_j)}. \quad (2)$$

Additionally, we would like the stability indicator to be independent of specific values of  $i$  and  $n$ . Hence, we normalize it by the expected Canberra distance between random rankings, denoted by  $\text{Ca}(n, i)$ . It can be approximated (Jurman et al., 2008) as

$$\text{Ca}(n, i) \approx \frac{(i+1)(2n-i)}{n} \log 4 + \frac{i(1+i)}{n} + 2i - 3, \quad (3)$$

which we make use of when  $i \geq 8$  and the computation of the exact value becomes intractable. For  $i \geq 8$ , the relative error of the approximation (3) is smaller than 1%. Our final stability indicator is thus the curve consisting of points calculated as

$$\left(i, \frac{St(\mathcal{S}^i)}{\text{Ca}(n, i)}\right), \quad (4)$$

for  $1 \leq i \leq n$ , that represent the relative change of distance between top- $i$  lists w.r.t. the expected top- $i$  distance.

## EMPIRICAL EVALUATION OF FFA/RFA CURVE PROPERTIES

We start with the experiments on synthetic datasets. In such laboratory conditions, one has full control over the data, can establish the ground truth feature ranking, and produce rankings of different quality. Such a setting will facilitate proper assessment of our proposed feature ranking evaluation method. Before proceeding to the experiments, we briefly describe the constructed synthetic datasets. The detailed description of the datasets is given in Appendix A1.

We construct three datasets named `single`, `pair` and `combined`. Each of them consists of 1000 instances and 100 features. The relevant features in `single` dataset are individually correlated to target, the relevant features in the `pair` dataset are related to the target via XOR relation, and the relevant features in the `combined` dataset are the union of the relevant features in the `single` and `pair` datasets. The rest of the features in the datasets are random noise.

### Evaluation by Randomising the Ground Truth Ranking

The appropriateness of the proposed method is first demonstrated on a family of feature rankings that contain more and more noise. By doing that, we can show that the lower and lower quality of feature rankings is reflected in the FFA and RFA curves, and thus detected by the method.

We start with the ground-truth ranking  $\mathbf{R}_{GT}$  and perturb it as follows. First, we select a proportion  $\theta$  of randomly selected features which are then assigned random relevances, drawn uniformly from the unit interval. The other features preserve their ground truth relevance. This results in a ranking  $\mathbf{R}_{\theta}$ .

### Experimental setup

We use the aforementioned `single`, `pair` and `combined` datasets. The following amounts  $\theta$  of noise are introduced into the ground truth ranking:  $\theta \in \{0.05, 0.1, 0.15, 0.2, 0.3, 0.5, 1\}$ . The value  $\theta = 1$  corresponds to a completely random ranking.

For every value of  $\theta$ , we estimate the expected values of the FFA/RFA curves that belong to the ranking  $\mathbf{R}_{\theta}$ , by first generating  $m = 100$  realizations of the ranking, and second, (point-wise) averaging of the error estimates of the obtained predictive models.

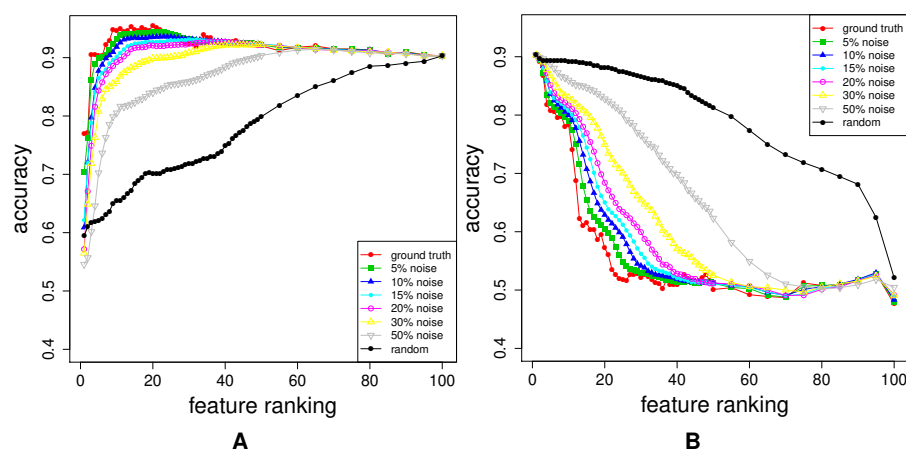
For constructing FFA/RFA curves, SVMs were used, as noted and justified at the end of the section *Analysis of Different Learning Methods to Construct FFA and RFA Curves*. The curves were constructed via 10-times stratified 10-fold cross validation, using different random seeds.

### Results

The obtained FFA and RFA curves are shown in Fig. 3 that gives the results for the dataset `combined`. The results for the datasets `single` and `pair` are similar. In addition to the curves that belong to the rankings  $\mathbf{R}_{\theta}$  with different amounts of noise, ground truth ranking is also shown.

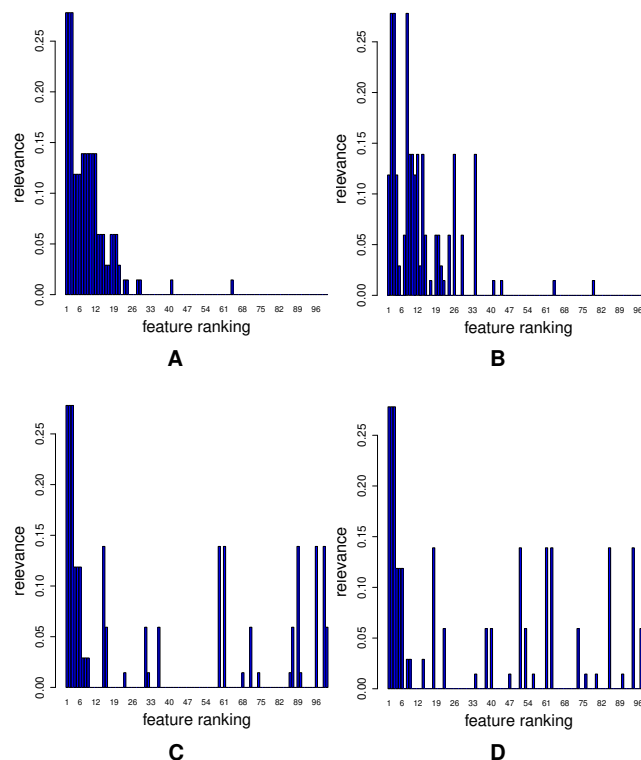
Both, FFA curves (Fig. 3A) and RFA curves (Fig. 3B) correctly detect different amount of noise  $\theta$ : the higher the  $\theta$ , the more distant are FFA and RFA curve of  $\mathbf{R}_{\theta}$  to the curves of ground truth ranking. The independent confirmation of these results are given in Appendix A2.

Additionally, note that FFA curves cannot give all the information about the ranking: Had we not plotted the RFA curves in Fig. 3B, we would not have a proof that all of the rankings misplace some relevant features (check the considerable decrease in accuracy just before the 100th feature).



**Figure 3.** Dataset `combined`: Forward feature addition curves (A), and reverse feature addition curves (B). The curves for the noisy rankings  $\mathbf{R}_{\theta}$  ( $0.05 \leq \theta \leq 1$ ) and the ground truth ranking are shown.





**Figure 4.** Distribution of relevant features for each of the four ranking methods: Relief F (A), Random Forests (B), Info Gain (C) and SVM-RFE (D).

## Analysis of Different Learning Methods to Construct FFA and RFA Curves

According to Algorithm 1, the error curve estimates depend not only on the feature ranking method, but also on the learning method used to construct the predictive models. In this section, we investigate which learning methods (learners) are suitable to construct the FFA and RFA curves. Note that we are not searching for a learner that would produce the most accurate predictive models. Rather, the requirement for the learner to be used in this context is that it should produce predictive models that exploit all the information that the features contain about the target concept, and can thus distinguish between feature rankings of different quality.

### Experimental Setup

When comparing the FFA and RFA curves of different ranking methods, constructed with different learners, we used the combined dataset described in detail in Appendix A1. We consider the following four feature ranking methods.

**Information gain**, where we calculate the information gain of each feature  $F_i$  as  $IG(F_i) = H(F_t) - H(F_t|F_i)$ .

**SVM-RFE** uses a series of support vector machines (SVMs) to recursively compute a feature ranking. A linear SVM was employed, as proposed by (Guyon et al., 2002). Following (Slavkov et al., 2018), we set  $\epsilon = 10^{-12}$  and  $c = 0.1$ .

**Relief F** algorithm as proposed by Robnik-Šikonja and Kononenko (Robnik-Šikonja and Kononenko, 2003). The number of neighbors was set to 10, and the number of iterations was set to 100%.

**Random Forests**, which can be used for estimating feature relevance as described by Breiman (Breiman, 2001). A forest of 100 trees was used, constructed by randomly choosing rounded up  $\log_2$  of the number of features.

We compare the above ranking methods by using different learners to produce classifiers and generate error estimates for the FFA and RFA curves. More specifically, we consider **Naïve Bayes** (John and Langley, 1995); **Decision Trees** (Quinlan, 1993); **Random Forests** (Breiman, 2001): the number of trees was set to 100, and in each split  $\log_2$  of the number of features are considered; **SVMs** (Cortes and Vapnik,

1995): a polynomial (quadratic) kernel was employed, with the  $\varepsilon = 10^{-12}$  and  $c = 0.1$ ; **k-NN** (Aha and Kibler, 1991) with a value of  $k = 10$ .

The curves were constructed via 10-times stratified 10-fold cross validation, using different random seeds. The obtained FFA and RFA curve comparisons of the four feature ranking methods obtained by each of the five learning methods, are presented in the following section.

# Results

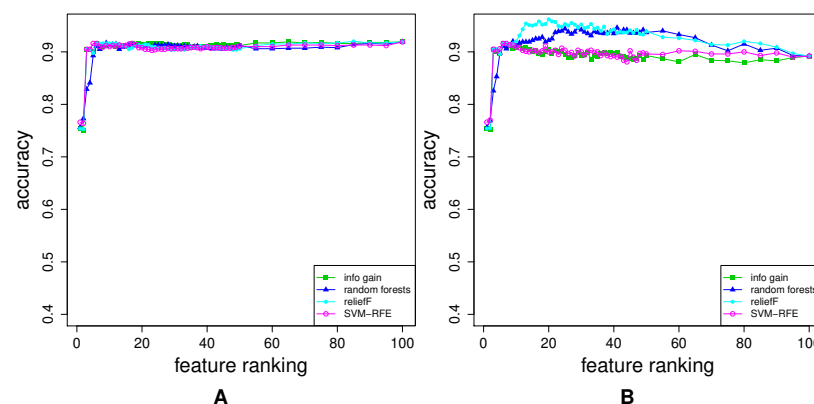
The rankings are shown in Fig. 4, where each graph represents the distribution of the ground truth relevance values. The y-axis depicts the ground truth relevance value (5). Each point,  $i$ , represents the  $i$ -th ranked feature,  $F_{(i)}$ , as determined by the feature ranking method.

We can see that the rankings fall into two groups: in Figs. 4A and 4B, highly relevant features are concentrated on the left, while in the Figs. 4C and 4D, they are evenly spread.

ReliefF and Random Forests (Figs. 4A and 4B) are thus clearly better than Info Gain and SVM-RFE (Figs. 4C and 4D). Hence, the FFA and RFA curves should at least differentiate between the two groups of the rankings. However, there should be visible difference also between Relief and Random Forests at the beginning of the ranking. The detailed comparative evaluation of the obtained feature rankings is given in Appendix A3.

In the case of FFA curves, the learners can be divided into two groups: FFA curves produced by Naïve Bayes, Decision Trees and Random Forests cannot capture any difference between rankings, whereas those produced by SVMs and k-NN can. It suffices to show one representative graph for each group (those for Naïve Bayes and k-NN are shown in Fig. 5), since there are no significant differences among the learning methods in the same group<sup>1</sup>. The FFA curves produced by these two learners have all the desirable properties: the curves for Relief and Random Forests are better than those of Info Gain and SVM-RFE. Additionally, at the beginning, the Random Forest curve is under the curve of Relief.

The reason why, for example, the Naïve Bayes classifier does not show any difference between rankings, is the fact that it can not use the information from the interactions of higher order. Namely, it assumes feature independence. Hence, it is not appropriate for use in the considered context.



**Figure 5.** Comparison of FFA curves for the four different ranking methods for the combined dataset. The curves were obtained by using the Naïve Bayes (A) and k-NN (B).

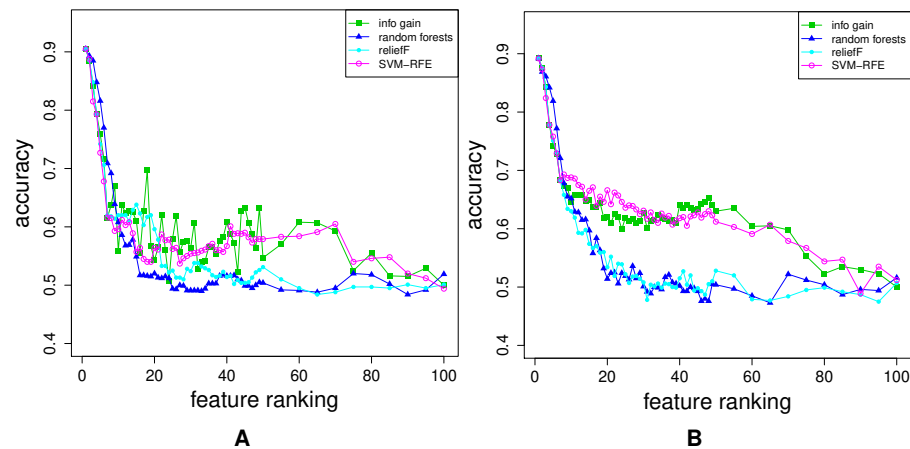
If we proceed to RFA curves, again, the Naïve Bayes classifier does not show any difference between rankings, whereas the other four methods do. We prefer Random Forests, SVMs and k-NN over Decision Trees in the case of RFA curves, because Decision Trees generate quite unstable curves, as shown in Fig. 6A. In Fig. 6B, the RFA curves of k-NN are shown. Again, there is no quantitative difference between them and the RFA curves generated by SVMs<sup>2</sup>.

Tu sum up, one can use

- SVMs and k-NN models, for constructing FFA curves,

<sup>1</sup>Compare, for example, Fig. 5B (obtained by k-NN) with Fig. A1A (obtained by SVMs), which is given in Appendix A3 (note that Fig. A1A also contains the random ranking curve).

<sup>2</sup>Compare Fig. 6B to Fig. A1B (given in Appendix A3).



**Figure 6.** Comparison of RFA curves for the four different ranking methods for the combined dataset. The curves were obtained by using Decision Trees (A) and k-NN (B).

- SVMs, k-NN and Random Forests, for constructing RFA curves.

Thus, only k-NN and SVMs are appropriate for constructing both FFA and RFA curves. Since one should typically use approximate k-NN when the number of features is extremely high (Muja and Lowe, 2009), we use SVMs (with the settings described here) as the learner for constructing the FFA/RFA curves in all the remaining experiments in this work.

### Discussion

We have to give some additional notes about choosing the best method, when, for example, different learning methods prioritize different rankings, which is possible since some learning method might make use of some features, whereas another learning method can make better use of some others.

If we have computed feature rankings to learn a classifier that uses only a subset of (top-ranked) features and we have already decided on which classifier to use, we should use the same (type of) classifier to construct the curves, because we want to use the features that the chosen learning method prioritizes.

Second, if our motivation for computing the feature rankings is to discover all relevant features for a given problem (e.g., the genes that influence the patients' clinical state), and learning method *A* prioritizes the ranking  $\mathbf{R} = (x_1, x_2, \dots)$  over the ranking  $\mathbf{R}' = (x'_1, x'_2, \dots)$ , whereas learning method *B* prioritizes  $\mathbf{R}'$  over  $\mathbf{R}$ , this means that  $x_1, x_2, x'_1$  and  $x'_2$  are important (provided that both learners achieve similar accuracy), so we can include them all in the subsequent experiments (and thus use both learning methods).

The decision about which among the two appropriate methods to use – k-NN or SVMs – might also depend on the properties of the dataset at hand. As mentioned before, k-NN could be too time-consuming when the number of features is extremely high. On the other hand, if the number of instances is high, SVMs could be too time-consuming, but speed-ups are possible (Tsang et al., 2005). As for the noise, both methods are quite robust (Wang et al., 2018; Xu et al., 2009), so this is not among the most influential factors.

## FEATURE RANKING COMPARISON ON REAL-WORLDS DATASETS

In this section, we move from the synthetic data and show the appropriateness of the proposed feature ranking evaluation method on the real-world data with unknown relevant and irrelevant features. To be consistent with the synthetic-data experiments, we evaluate the same four feature ranking methods as before, and compare them to the random feature rankings which now serve as the only baseline.

### Datasets Description

In this extensive study, 35 classification benchmark problems are used. They come in two bigger groups. The first group has been part of the experiments in (Slavkov et al., 2018) and except for *aapc* (Džeroski et al., 1997), *water* and *diversity* (Džeroski et al., 2000), mostly originates from the UCI data repository (Newman and Merz, 1998). This benchmark problems have higher number of instances (up

to 5000) and not extremely high number of features (up to 280). These problems cover various domains: health, ecology, banking etc.

The second group is newly included and contains 11 high-dimensional micro-array benchmark problems (Mramor et al., 2007) (up to 12625 features) with lower number of examples (up to 110). The main properties of the data are given in Tab. 1.

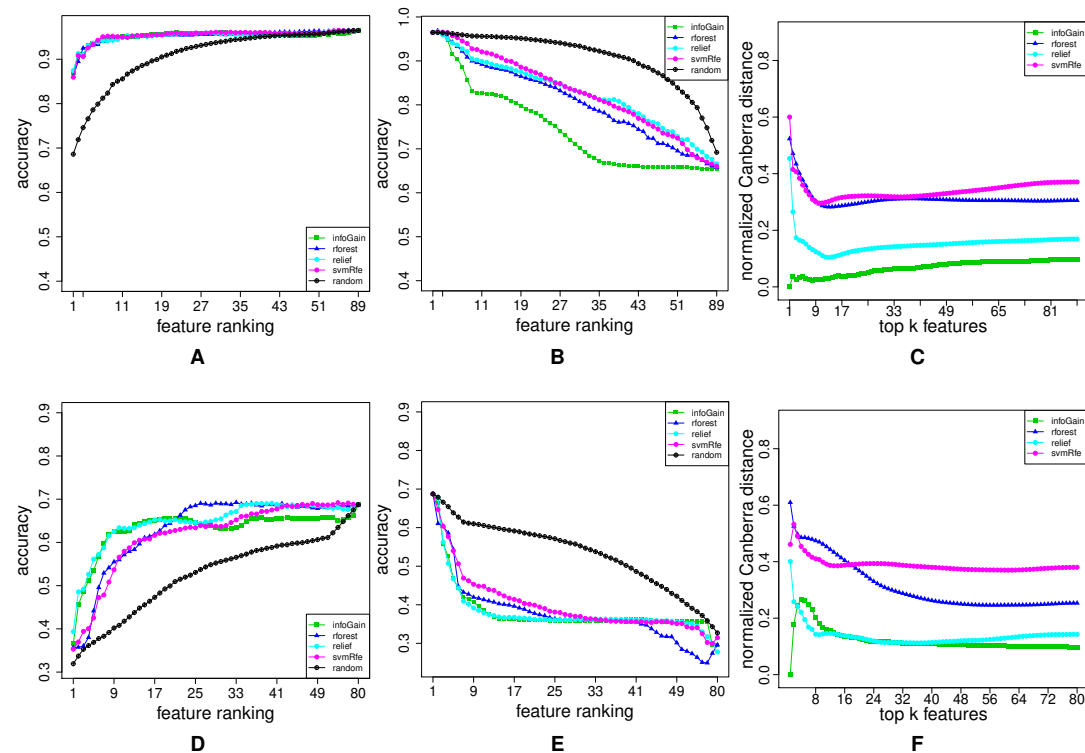
**Table 1.** Properties of the benchmark datasets: number of instances, number of features, number of discrete/numeric attributes, and number of different class values. The datasets with a considerably high number of features are listed under the dashed line.

Dataset	#Inst.	#Feat.	(D/N)	#Cl.
aapc	335	84	(83/1)	3
arrhythmia	452	280	(73/206)	16
australian	690	14	(8/6)	2
balance	625	4	(0/4)	3
breast-cancer	286	9	(9/0)	2
breast-w	699	9	(9/0)	2
car	1728	6	(6/0)	4
chess	3196	36	(36/0)	2
diabetes	768	8	(0/8)	2
diversity	292	86	(0/86)	5
german	1000	20	(13/7)	2
heart	270	13	(6/7)	2
heart-c	303	13	(7/6)	5
heart-h	294	13	(7/6)	5
hepatitis	155	19	(13/6)	2
image	2310	19	(0/19)	7
ionosphere	351	34	(0/34)	2
iris	150	4	(0/4)	3
sonar	208	60	(0/60)	2
tic-tac-toe	958	9	(9/0)	2
vote	435	16	(16/0)	2
water	292	80	(0/80)	5
waveform	5000	21	(0/21)	3
wine	178	13	(0/13)	3
amlPrognosis	54	12625	(0/12625)	2
bladderCancer	40	5724	(0/5724)	3
breastCancer	24	12625	(0/12625)	2
childhoodAll	110	8280	(0/8280)	2
cmlTreatment	28	12625	(0/12625)	2
colon	62	2000	(0/2000)	2
dlbcl	77	7070	(0/7070)	2
leukemia	72	5147	(0/5147)	2
ml	72	12533	(0/12533)	3
prostate	102	12533	(0/12533)	2
srbc	83	2308	(0/2308)	4

## Experimental Setup

We construct the curves that base on the feature ranking methods described in experimental setup part of Sec. *Analysis of Different Learning Methods to Construct FFA and RFA Curves*, and the curves that belong to the completely random ranking (i.e., expected curves) which serve as a baseline. For the actual construction of the curves (once the ranking is obtained), support vector machines were used, as described and justified in the results in Section *Analysis of Different Learning Methods to Construct FFA and RFA Curves*. The curves were constructed via 10-times stratified 10-fold cross validation, using different random seeds.

The expected error curves for random rankings were produced by generating 100 random rankings for each dataset under consideration. For each random ranking, error curves were produced and the average of the error values was used as the expected error. This was done in a manner similar to the one described in Sec. *Evaluation by Randomising the Ground Truth Ranking*. As mentioned in Sec. *The Evaluation Method*, building FFA/RFA curves by adding the features one by one to large feature subsets  $S^i$  and  $S_i$ , might be too costly when  $n$  is big enough. In this set of experiments, we use the following procedure. We add  $\delta(i)$  features to the subset, where  $\delta(i)$  is defined as follows:  $\delta(i) = 1$  if  $1 \leq i \leq 50$ ,  $\delta(i) = 5$  if  $50 < i \leq 500$ , and  $\delta(i) = n//20$  otherwise, where  $//$  denotes integer division.



**Figure 7.** Ranking quality assessment for datasets *breast-w* (first row) and *water* (second row) in terms of the FFA (first column) and RFA curves (second column), and rankings' stability estimates (third column). The FFA/RFA curves are obtained by using SVMs.

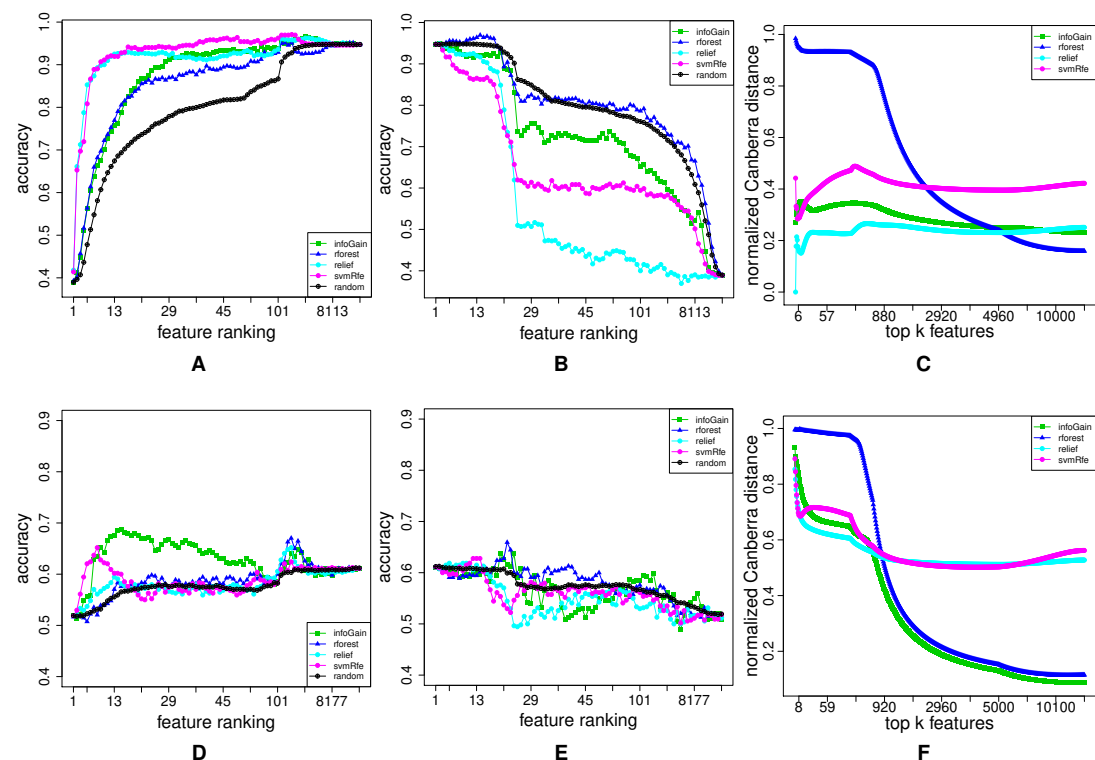
## Results

In this section, we show representative examples of three types of curves: FFA, RFA and stability curves. The curves are shown for two datasets with lower and two with higher number of features. The graphs for the other datasets can be found in Appendix A4 in Figs. A2, A3, A4, A5, A6, A7, A8, A9, A10, A11, and A12.

We start with the *breast-w* dataset. The FFA/RFA curves in Figs. 7A and 7B show that both types of curves are needed in order to evaluate the ranking completely. The FFA-curves suggest that all feature ranking algorithms (except for the random ranking) place only the relevant features at the beginning, since there is practically no difference if we compare the accuracy of the 89-feature (all) model and, e.g., the 11-feature model. However, the RFA-curves show that all feature ranking algorithms - except for Info Gain - misplace some relevant features, since the Info-Gain-ranking-based models have the lowest accuracy by far in the RFA curves. Also, in the case of Info Gain, the accuracy cease to decrease after only cca. 40 top-ranked features were removed.

Fig. 7C shows that Info Gain produces also the most stable rankings. We can see that the top-ranked feature is always the same, since the stability index of the Info Gain equals 0 at the point  $k = 1$ . The second most stable algorithm is ReliefF, the third is Random Forest and the least stable is SVM-RFE, but the difference between Random Forest and SVM-RFE is not considerable.

Let us now take a look at the curves for the dataset *water*. From the FFA-curves in Fig. 7D, we see that ReliefF, Info Gain and Random Forest ought to have the same 21 top-ranked features, and as a consequence, the same last  $59 = 80 - 21$  features too. However, the first 21 features are ordered better by Info Gain and ReliefF, while the last 59 are more properly ordered by Random Forest. We can conclude that none of the rankings is ideal, but we can come close to the ideal one (in terms of FFA-curves), if we combine the first part of the ReliefF (or Info Gain) and the second part of Random Forest. This claim is also confirmed by the RFA-curves of Info Gain and ReliefF (Fig. 7E): these two algorithms indeed misplace some relevant features, since the accuracy of the model abruptly decreases and the end of the ranking.



**Figure 8.** Ranking quality assessment for datasets *mll* (first row) and *amlPrognosis* (second row) in terms of the FFA (first column) and RFA curves (second column), and rankings' stability estimates (third column). The FFA/RFA curves are obtained by using SVMs.

Fig. 7F suggests that we should prefer Info Gain and ReliefF over Random Forest since they are more stable. However, we can also notice that Random Forest is the least stable at the beginning of the ranking but its stability increases when the number of features gets larger.

We begin the analysis of the high-dimensional datasets with *mll*. Fig. 8B shows that Random Forest completely misplaces some relevant features, since its RFA-curve mostly goes above the random-ranking one. Even though it is evident from Random Forest's FFA-curves that some relevant features are also successfully captured, Random Forest produces the worst ranking. Info Gain is slightly better, whereas ReliefF and SVM-RFE are again the best algorithms. From the FFA-curves, we can conclude that SVM-RFE places more features with higher relevance at the beginning of the ranking (its curve is higher than ReliefF's), while RFA-curves reveal that SVM-RFE also misses some quite relevant features: ReliefF's curve is far below SVM-RFE's. Fig. 8C shows that ReliefF is considerably more stable than SVM-RFE, hence we prefer the former over the latter on the *mll* dataset.

The last example will show that sometimes, the understanding of the results is not that easy. In the Figs. 8D and 8E, the FFA and RFA curves for the *amlPrognosis* dataset are presented. In this case, only Info Gain performs considerably better than random ranking in terms of FFA-curves. SVM-RFE is also able to find some relevant features at the beginning (peak of its curve at 10 features), but after that, the models' accuracy decreases, hence mostly noisy features are positioned here. Some relevant features are again placed by ReliefF, Info Gain and Random Forest also around the 2000th place (local peak of their curves in the right part of the FFA-curve). RFA-curves confirm that there is indeed much noise in these data, since removing features does not result in an (at least approximately) decreasing curve.

It may not come as a surprise that all ranking algorithms produce rankings that are very unstable at the beginning (Fig. 8F), but it is interesting that after approximately 1000 features, Info Gain and Random Forest produce quite stable rankings even though they have low quality. The reason for both low quality of rankings and their instability is probably the low number of instances accompanied by a high number of features (54 and 12625 respectively).

## CONCLUSIONS

We have proposed a method for evaluating and comparing feature rankings. The method is based on constructing two chains of predictive models that are built on two families of nested subsets of features. The first family of subsets are the sets of top-ranked features, while the second family consists of sets of bottom-ranked features. The predictive performances of the models form a forward feature addition (FFA) curve in the former case, and reverse feature addition (RFA) curve in the latter case.

We show in our experiments that both types of curves are necessary when comparing the rankings: FFA curves detect whether important features are placed at the beginning of the ranking, whereas RFA curves detect whether important features can still be found at the end of the ranking.

In the set of experiments, we show the usefulness of the proposed evaluation method and its sensitivity to the rankings of different quality on synthetic data. The second set of experiments shows which of the learning methods are appropriate for building the FFA and RFA curves (SVMs, k-NN) and which are not (Naïve Bayes, Decision Tree, Random Forest). In the third set of experiments on synthetic data, we test several feature ranking algorithms and examine their properties. Considering data with different properties, we show that ReliefF algorithm outperforms the other investigated approaches, both in terms of detecting relevant features and in terms of stability of the feature rankings it produces.

Moreover, we show the usefulness of the proposed approach in real-world scenarios. We evaluate feature rankings computed by four feature ranking algorithms on 35 classification benchmark problems. The results reveal that there is no feature ranking algorithm that would dominate the others on every dataset.

A possible disadvantage of the proposed method is that it can be computationally quite intensive, if we want to construct the curves in full resolution. Namely, every point of a FFA or RFA curve comes at the cost of building and evaluating a predictive model. However, as justified in the method description, the full resolution is, especially when the number of features is really high, not necessary, and, moreover, the construction of the curve can be also easily parallelized.

The work presented in this paper can continue in many directions. First, of all, the proposed methodology could use other error measures, since accuracy is appropriate only for the task of classification when the distribution of target variable is approximately uniform. The strong modularity of the FFA/RFA curves allows for their use in any other predictive modeling task, e.g., for the task of regression, we could use root mean squared error instead of accuracy. However, even though there exists a regression version of most of the learners, which are considered for constructing the curves, experiments should be repeated on those cases, since the conclusions about, for example, the most appropriate learner for constructing the curves, can be different. Moreover, the method can be adapted not only to the regression setting, but also to different tasks of structured output prediction (Bakır et al., 2007) and time series prediction.

## REFERENCES

- Aha, D. and Kibler, D. (1991). Instance-based learning algorithms. *Machine Learning*, 6:37–66.
- Arceo-Vilas, A., Fernandez-Lozano, C., Pita, S., Pérttega-Díaz, S., and Pazos, A. (2020). A redundancy-removing feature selection algorithm for nominal data. *PeerJ Computer Science*.
- Bakır, G. H., Hofmann, T., Schölkopf, B., Smola, A. J., Taskar, B., and Vishwanathan, S. V. N., editors (2007). *Predicting Structured Data*. The MIT Press, Cambridge, Massachusetts.
- Biesiada, J., Duch, W., Kachel, A., and Paucha, S. (2005). Feature ranking methods based on information entropy with parzen windows. In *International Conference on Research in Electrotechnology and Applied Informatics*, Katowice, Poland.
- Boucheham, A. and Batouche, M. (2014). Robust biomarker discovery for cancer diagnosis based on meta-ensemble feature selection. In *2014 Science and Information Conference*, pages 452–560.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45:5–32.
- Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3):273–297.
- Duch, W., Wiecek, T., and Biesiada, J. (2004). Comparison of feature ranking methods based on information entropy. In Fawcett, T. and Mishra, N., editors, *IEEE International Conference on Neural Networks - Conference Proceedings*, volume 2, pages 1415–1419. IEEE.
- Džeroski, S., Demšar, D., and Grbović, J. (2000). Predicting chemical parameters of river water quality from bioindicator data. *Applied Intelligence*, 13:7–17.
- Džeroski, S., Potamias, G., Moustakis, V., and Charissis, G. (1997). Automated revision of expert rules



- 487 for treating acute abdominal pain in children. In *Artificial intelligence in medicine - AIME, LNCS 1211*,  
488 pages 98–109.
- 489 Furlanello, C., Serafini, M., Merler, S., and Jurman, G. (2003). Entropy-based gene ranking without  
490 selection bias for the predictive classification of microarray data. *BMC Bioinformatics*, 4:54.
- 491 Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. (2002). Gene selection for cancer classification using  
492 support vector machines. *Machine Learning*, 46:389–422.
- 493 Guzmán-Martínez, R. and Alaiz-Rodríguez, R. (2011). Feature selection stability assessment based on  
494 the Jensen-Shannon divergence. *Lecture Notes in Computer Science*, 6911:597–612.
- 495 Henzgen, S. and Hüllermeier, E. (2015). Weighted rank correlation: A flexible approach based on fuzzy  
496 order relations. In Appice, A., Rodrigues, P. P., Santos Costa, V., Gama, J., Jorge, A., and Soares,  
497 C., editors, *Machine Learning and Knowledge Discovery in Databases*, pages 422–437. Springer  
498 International Publishing.
- 499 John, G. H. and Langley, P. (1995). Estimating continuous distributions in Bayesian classifiers. In *proc.*  
500 *Eleventh Conference on Uncertainty in Artificial Intelligence*, pages 338–345, San Mateo, CA. Morgan  
501 Kaufmann.
- 502 Jong, K., Mary, J., Cornuéjols, A., Marchiori, E., and Sebag, M. (2004). Ensemble feature ranking. In  
503 *PKDD - LNCS 2302*, pages 267–278.
- 504 Jurman, G., Merler, S., Barla, A., Paoli, S., Galea, A., and Furlanello, C. (2008). Algebraic stability  
505 indicators for ranked lists in molecular profiling. *Bioinformatics*, 24:258–264.
- 506 Kalousis, A., Prados, J., and Hilario, M. (2007). Stability of feature selection algorithms: A study on  
507 high-dimensional spaces. *Knowledge and Information Systems*, 12:95–116.
- 508 Khoshgoftaar, T. M., Fazelpour, A., Wang, H., and Wald, R. (2013). A survey of stability analysis of  
509 feature subset selection techniques. In *IEEE 14th International Conference on Information Reuse*  
510 *Integration (IRI)*, pages 424–431.
- 511 Lance, G. N. and Williams, W. T. (1966). Computer programs for hierarchical polythetic classification  
512 ('similarity analyses'). *The Computer Journal*, 9:60–64.
- 513 Lance, G. N. and Williams, W. T. (1967). Mixed-data classificatory programs i - agglomerative systems.  
514 *Australian Computer Journal*, 1.
- 515 Li, Z. and Gu, W. (2015). A redundancy-removing feature selection algorithm for nominal data. *PeerJ*  
516 *Computer Science*.
- 517 Liang, J., Yang, S., and Winstanley, A. (2008). Invariant optimal feature selection: A distance discriminant  
518 and feature ranking based solution. *Pattern Recognition*, 41:1429–1439.
- 519 Liu, T., Liu, S., Chen, Z., and Ma, W.-Y. (2003). An evaluation on feature selection for text clustering. In  
520 Fawcett, T. and Mishra, N., editors, *ICML*, pages 488–495. The AAAI Press, Menlo Park, California.
- 521 Mramor, M., Leban, G., Demšar, J., and Zupan, B. (2007). Visualization-based cancer microarray data  
522 classification analysis. *Bioinformatics*, 23(16):2147.
- 523 Muja, M. and Lowe, D. G. (2009). Fast approximate nearest neighbors with automatic algorithm  
524 configuration. In Ranchordas, A. and Araújo, H., editors, *VISAPP (I)*, pages 331–340. INSTICC Press.
- 525 Nardone, D., Ciaramella, A., and Staiano, A. (2019). A redundancy-removing feature selection algorithm  
526 for nominal data. *PeerJ Computer Science*.
- 527 Newman, C. B. D. and Merz, C. (1998). UCI repository of machine learning databases.  
528 <https://archive.ics.uci.edu/ml/datasets.html>. Accessed on: 2015-12-13.
- 529 Nilsson, R., Peña, J. M., Björkegren, J., and Tegnér, J. (2007). Consistent feature selection for pattern  
530 recognition in polynomial time. *Journal of Machine Learning Research*, 8:589–612.
- 531 Nogueira, S., Sechidis, K., and Brown, G. (2017). On the stability of feature selection algorithms. *Journal*  
532 *of Machine Learning Research*, 18(1):6345–6398.
- 533 Paoli, S., Jurman, G., Albanese, D., Merler, S., and Furlanello, C. (2005). Semisupervised profiling  
534 of gene expressions and clinical data. In *Proc. Sixth International Conference on Fuzzy Logic and*  
535 *Applications*, pages 284–289.
- 536 Quinlan, R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers, San Mateo,  
537 CA.
- 538 Robnik-Šikonja, M. and Kononenko, I. (2003). Theoretical and empirical analysis of ReliefF and RReliefF.  
539 *Machine Learning*, 53:23–69.
- 540 Saeys, Y., Abeel, T., and de Peer, Y. V. (2008). Robust feature selection using ensemble feature selection  
541 techniques. In *ECML/PKDD, LNCS 5212*, pages 313–325.



- Slavkov, I., Petković, M., Koccev, D., and Džeroski, S. (2018). Quantitative score for assessing the quality of feature rankings. *Informatica*, 42(1):43–52.
- Tsang, I. W., Kwok, J. T., and Cheung, P.-M. (2005). Core vector machines: Fast svm training on very large data sets. *Journal of Machine Learning Research*, 6:363–392.
- Verikas, A., Gelzinis, A., and Bacauskiene, M. (2011). Mining data with random forests: A survey and results of new tests. *Pattern Recognition*, 44:330–349.
- Wang, Y., Jha, S., and Chaudhuri, K. (2018). Analyzing the robustness of nearest neighbors to adversarial examples. In Dy, J. G. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5120–5129. PMLR.
- Xu, H., Caramanis, C., and Mannor, S. (2009). Robustness and regularization of support vector machines. *Journal of Machine Learning Research*, 10:1485–1510.

## A1

In this section, we explain how we generate our synthetic datasets. For simplicity, we take both the features  $F_i$  and the target  $F_t$  to be binary and take values from the set  $\{0, 1\}$ . We then partition the set of features  $\mathcal{F}$  into feature interaction subsets  $\mathcal{F}_{\text{int}}$  of cardinality one and two. The feature sets with cardinality one are single features  $F_i$  that are in an individual interaction with the target  $F_t$ , while the features from the interaction sets with cardinality two determine the value of the target by the *XOR* relation.

The examples are generated as follows. For each example, we first randomly (from a uniform distribution) set the value of the target feature  $F_t$ . After that, if  $\mathcal{F}_{\text{int}} = \{F_i\}$ , then the value of feature  $F_i$  is randomly chosen, so that  $P(F_i = F_t) = p$ . Otherwise, we have  $\mathcal{F}_{\text{int}} = \{F_i, F_j\}$ , and the values of the features  $F_i$  and  $F_j$  are randomly chosen, so that  $P(\text{XOR}(F_i, F_j) = F_t) = p$ .

We consider the probability values  $p \in \{0.8, 0.7, 0.6, 0.5\}$ . The feature sets with  $p = 0.5$  are in fact independent of the target  $F_t$ , and can be considered as irrelevant features.

With combinations of these feature interaction sets, three datasets were generated, each of them consisting of 1000 instances and 100 features in total.

The first dataset (named `single`) comprises only individually correlated features. The second dataset (named `pair`) contains relevant features related to the target via the *XOR* relation, as well as irrelevant features. The third (named `combined`) is a combination of the first two. It contains *XOR*-related features and individually correlated features.

In order to simulate the redundancy of features, which occurs in real datasets, the three datasets are constructed in the following way: If the set  $\mathcal{F}_{\text{int}}$  of relevant features is included in the dataset, we also include two redundant copies of  $\mathcal{F}_{\text{int}}$  in the dataset. Irrelevant features are realized independently of each other.

The properties of the generated datasets are summarized in Table 2, from which we can observe that there are 9 relevant features in the `single` dataset, 18 in the `pair` dataset, and 27 in the `combined` dataset.

**Table 2.** Properties of the synthetic datasets. If  $p > 0.5$ , #copies denotes the number of copies of the interaction set. In the last row where  $p = 0.5$ , #copies corresponds to the number of independently realised irrelevant features.

$ \mathcal{F}_{\text{int}} $	$p$	# copies in single	# copies in pair	# copies in combined
1	0.8	3		3
	0.7	3		3
	0.6	3		3
2	0.8		3	3
	0.7		3	3
	0.6		3	3
1	0.5	91	82	73

The ground-truth feature relevances  $\text{rel}(F_i)$  of the features  $F_i$  are defined as follows. First, the relevance of each feature group  $\mathcal{F}_{\text{int}}$  is defined as the mutual information between the group and  $F_i$ , namely  $\text{rel}(\mathcal{F}_{\text{int}}) = I(\mathcal{F}_{\text{int}}; F_i)$ . Second, for  $F_i \in \mathcal{F}_{\text{int}}$ , feature importances are obtained as

$$\text{rel}(F_i) = \text{rel}(\mathcal{F}_{\text{int}}) / |\mathcal{F}_{\text{int}}|. \quad (5)$$

For the particular three datasets, this ground-truth ranking  $R_{GT}$  should also result in the optimal FFA and RFA curves, but this may not be the case in general. In the next section, we give the results of comparing it to the other feature rankings.

## A2

When discussing the results in the subsection *Evaluation by Randomising the Ground Truth Ranking*, we showed that, when the level of noise  $\theta$  in the ranking increases, then i) the quality of the ranking  $R_\theta$  decreases, and ii) the rankings  $R_{GT}$  and  $R_\theta$  become more and more distant. However, for the second point, we need to define a distance  $\text{dist}(R_{GT}, R_\theta)$  between a noisy and the ground truth ranking. In the definition of  $\text{dist}(R_{GT}, R_\theta)$  we use the average Spearman rank correlation coefficient  $\rho(R_A, R_B)$  which is calculated as

$$\frac{1}{n-1} \sum_{i=1}^n \frac{(\text{rank}_A(F_i) - \overline{\text{rank}_A})(\text{rank}_B(F_i) - \overline{\text{rank}_B})}{\sigma_A \sigma_B}$$

where  $n$  is the number of features, therefore the average ranks  $\overline{\text{rank}_A}$  and  $\overline{\text{rank}_B}$  equal  $(n+1)/2$ . Standard variations  $\sigma_{A,B}$  are computed as  $\sigma_{A,B} = \sqrt{\sum_{i=1}^n (\text{rank}_{A,B}(F_i) - \overline{\text{rank}_{A,B}})^2 / (n-1)}$ . For a given  $\theta$ , the distance between rankings  $R_{GT}$  and  $R_\theta$  is then computed as

$$\text{dist}(R_{GT}, R_\theta) = 1 - \frac{1}{m} \sum_{t=1}^m \rho(R_{GT}, R_{\theta,t}) \quad (6)$$

where  $m$  is the number of different realizations  $R_{\theta,t}$  of the noisy ranking  $R_\theta$ .

Tab. 3 lists the values of the distances between the ground truth ranking  $R_{GT}$  and its noisy versions  $R_\theta$ . Note that, for all three synthetic datasets, the distances indeed increase when the amount of noise  $\theta$  increases.

**Table 3.** The distances  $\text{dist}(R_{GT}, R_\theta)$  for different  $\theta$  values, for each of the three synthetic datasets.

$\theta$	0.05	0.1	0.15	0.2	0.3	0.5	1
single	0.219	0.34	0.449	0.51	0.628	0.81	1.056
pair	0.126	0.239	0.327	0.397	0.519	0.726	1.091
combined	0.1	0.171	0.252	0.32	0.432	0.652	1.048

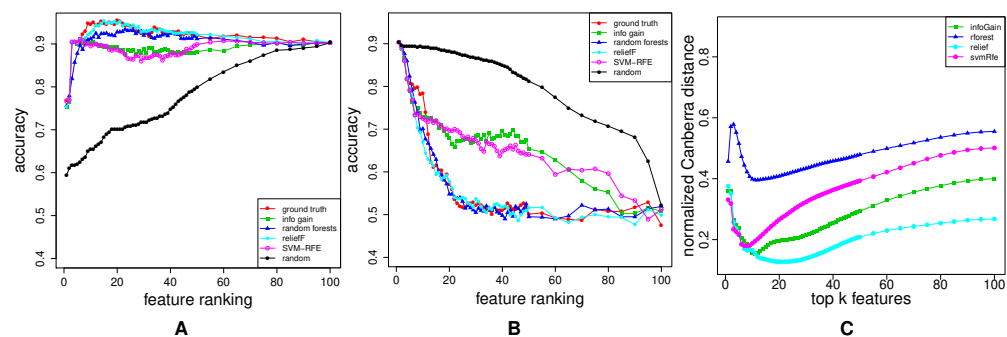
## A3

In Sec. *Evaluation by Randomising the Ground Truth Ranking*, we analyzed rankings of different quality (with different amounts of noise) by comparing them to the ground truth ranking. In a real-world setting, the ground truth ranking is unknown and the feature rankings are induced directly from data. Therefore, in this section we analyze feature rankings, produced by the four feature ranking methods from Sec. *Analysis of Different Learning Methods to Construct FFA and RFA Curves*, and the synthetic data described in Sec. *Generating Synthetic Data*.

When comparing the rankings, stability should also be taken into account as discussed earlier. Therefore, the stability indicator (4) is also included in the analysis.

### Experimental Setup

We have used the same parameter settings for the feature ranking algorithms as in Sec. *Analysis of Different Learning Methods to Construct FFA and RFA Curves*. As noted and justified in the corresponding Sec. *Results*, SVMs were used for constructing the FFA/RFA curves. The curves were constructed via 10-times stratified 10-fold cross validation, using different random seeds. To complement them, we also estimate the stability of each feature ranking algorithm by using the stability indicator described in Sec. *Stability of feature ranking*. All feature ranking methods were tested only on the combined dataset.



**Figure A1.** Ranking quality assessment for dataset combined in terms of the FFA (A) and RFA curves (B), and rankings' stability estimates (C). The FFA/RFA curves are obtained by using SVMs.

## Results

For our analysis, we consider three types of graphs. The first two types are FFA curves (Fig. A1A) and RFA curves (Fig. A1B). The third is the stability estimate graph (Fig. A1C) where the y-axis refers to the value of the stability indicator (4): the higher the value, the less stable the ranking method. Each point,  $k$ , on the x-axis represents the size of the considered feature subsets, consisting of the top ranked  $k$  features.

Upon a visual inspection of the overall results in Fig. 4, we can conclude that all of the feature ranking methods can correctly detect the features individually related to the target. However, Info Gain and SVM-RFE (Figs. 4C and 4D, respectively) exhibit random behavior for the XOR features, i.e., are unable to assign proper relevance values to them. Random Forests (Fig. 4B) separate relevant from irrelevant features, but the ordering of the relevant features is mixed. Finally, ReliefF (Fig. 4A) provides the ranking that is the closest to the ground truth.

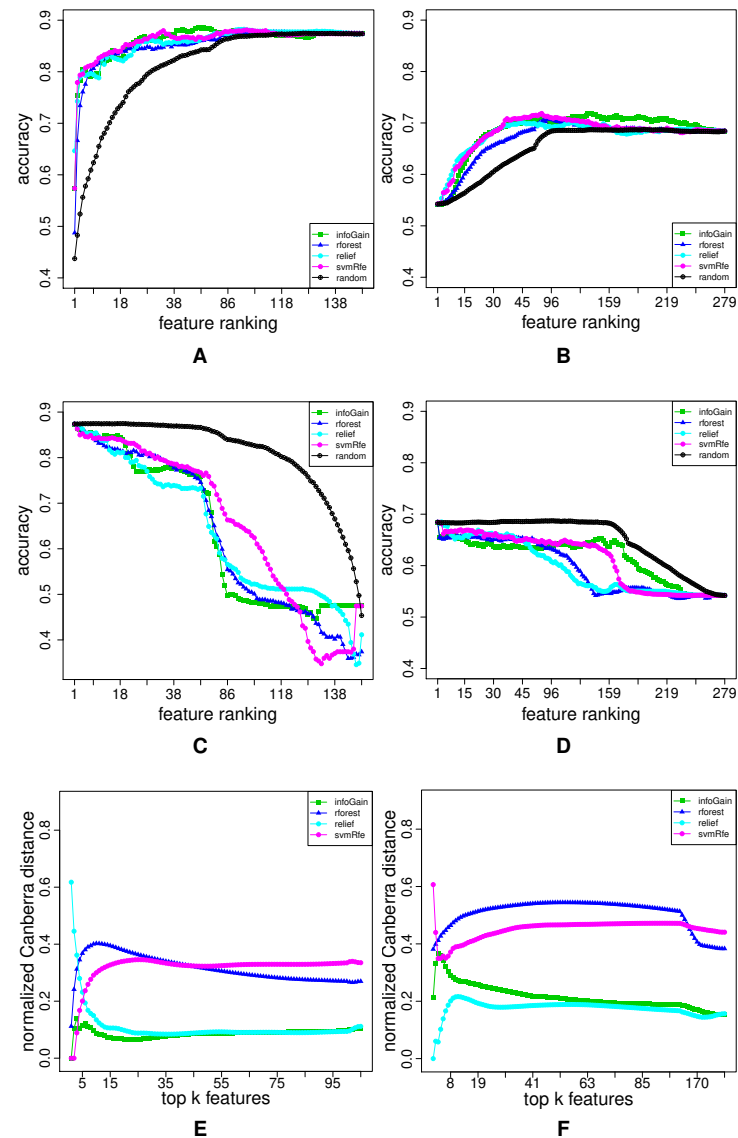
These differences in behavior among the different ranking methods are also clearly reflected in the FFA and RFA curves in Figs. A1A and A1B. In Fig. A1B, the RFA curves for Info Gain and SVM-RFE have a similar behavior: Namely, a linearly increasing accuracy (from right to left) in the region where the relevant features are randomly distributed and a sharp increase in accuracy in the region where the individually relevant features are included. On the other hand, the RFA curves of both random forests and ReliefF remain first constant and then increase abruptly when the top-ranked features are included. These two groups of methods can be also distinguished from the FFA curves. The FFA curves of all methods are first increasing abruptly and then slightly decreasing but the FFA curves of ReliefF and random forests increase during more steps and reach higher accuracy than Info Gain and SVM-RFE. This clearly indicates that while Info Gain and SVM manage to identify a proportion of the relevant features and put them at the top of the ranking, this proportion is nevertheless smaller than the one identified by random forests and ReliefF.

FFA and RFA curves undoubtedly allow us to compare the quality of the different ranking algorithms. The FFA/RFA curves of all methods are clearly better than the curves of the random ranking. The ReliefF ranking algorithm, however, clearly outperforms the other methods. It has the best error curves, i.e., the curves that are the closest to the ground truth ranking. The second best method are random forests: they exhibit very similar performances, but show a slightly worse FFA curve. Both Info Gain and SVM-RFE are clearly inferior in terms of both FFA and RFA curves.

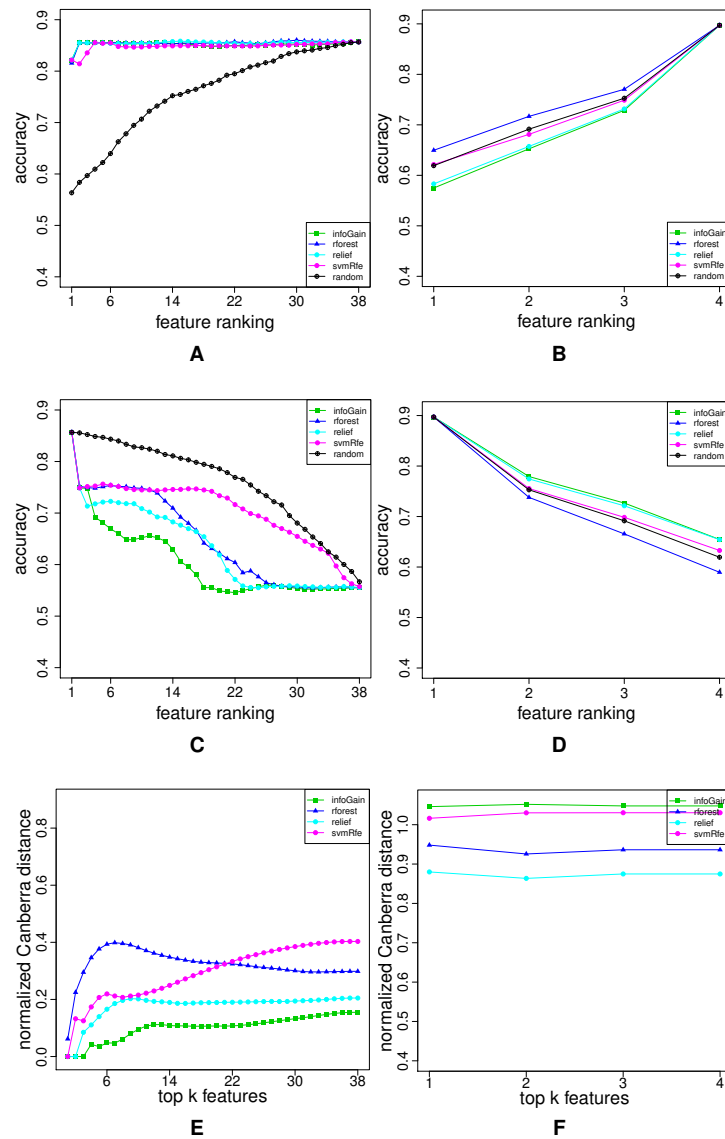
Stability-wise, as seen in Fig. A1C, all of the algorithms are stable in the region of the relevant features that they can detect, except for random forests, which has an instability peak exactly in this region. This means that random forests are in this case capable of detecting all the relevant features, but are highly unstable in the estimation of their ordering. Further inspection reveals that Relief generates not only the best rankings in terms of FFA/RFA curves, but also the most stable ones.

## A4

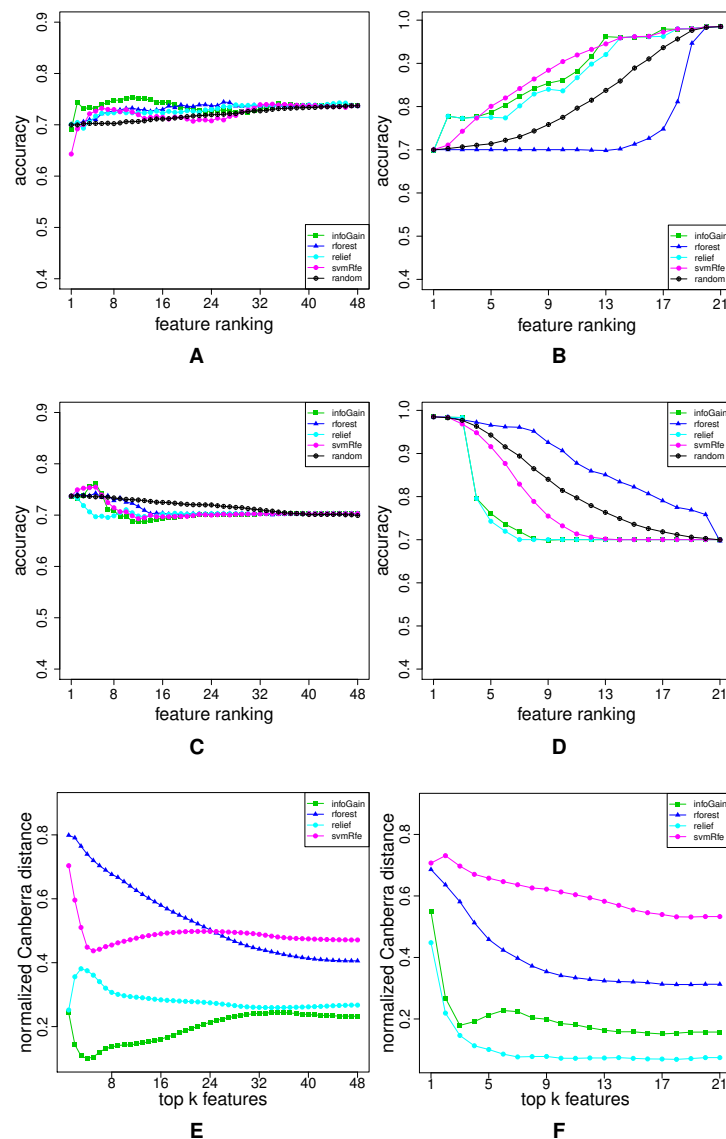
In this section, we provide detailed per dataset results from the experimental study.



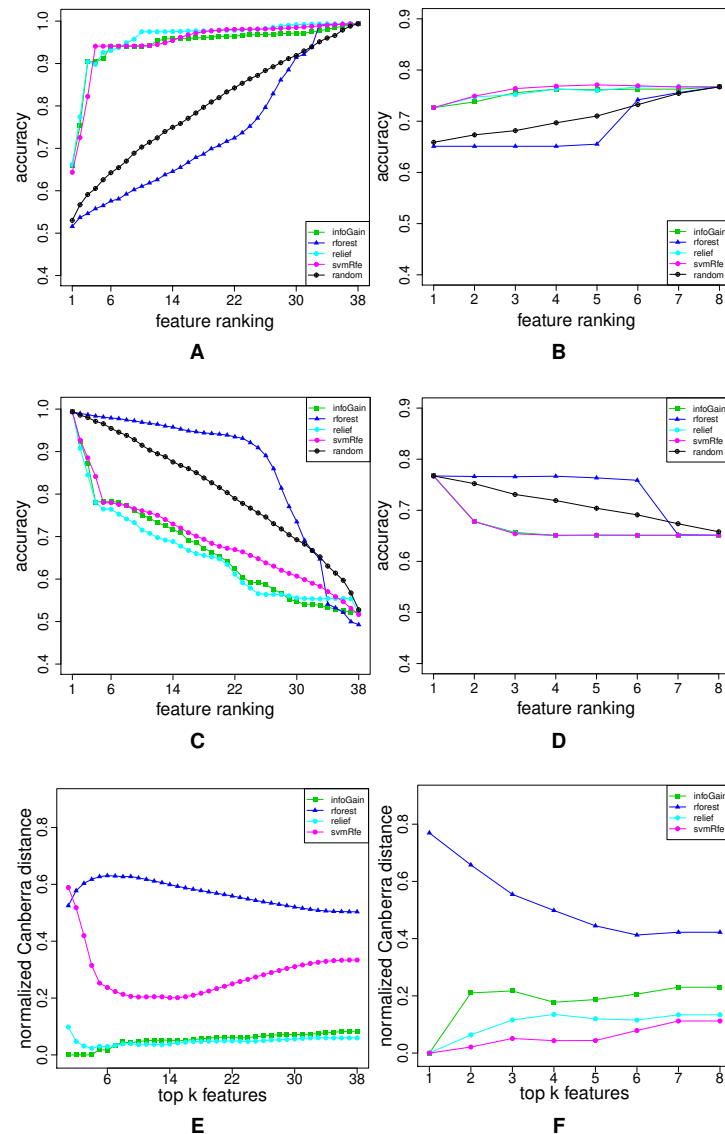
**Figure A2.** Ranking quality assessment for datasets *aapc* (first column) and *arrhythmia* (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



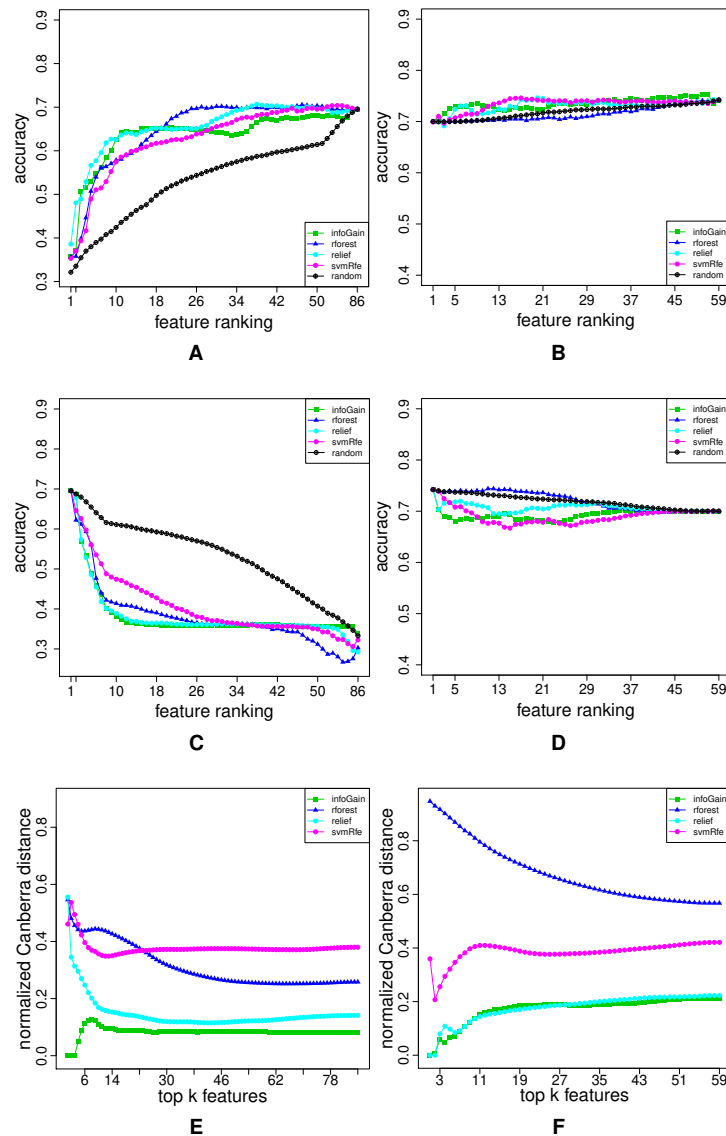
**Figure A3.** Ranking quality assessment for datasets australian (first column) and balance (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



**Figure A4.** Ranking quality assessment for datasets `breast-cancer` (first column) and `car` (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.

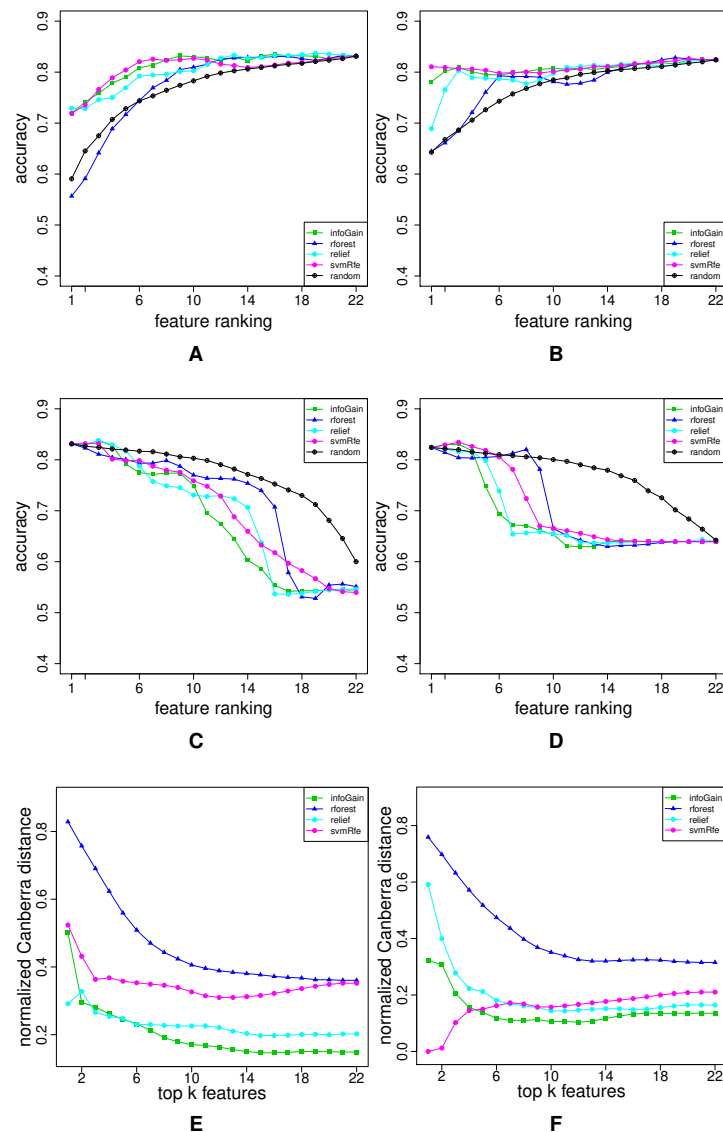


**Figure A5.** Ranking quality assessment for datasets *chess* (first column) and *diabetes* (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.

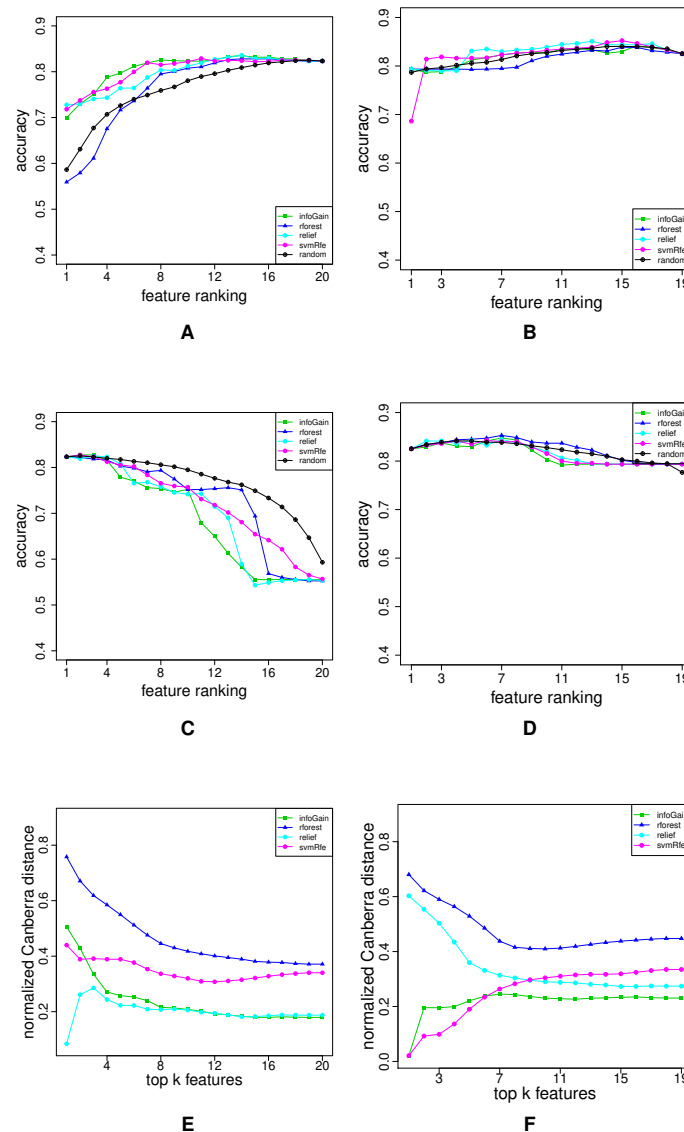


**Figure A6.** Ranking quality assessment for datasets `diversity` (first column) and `german` (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.

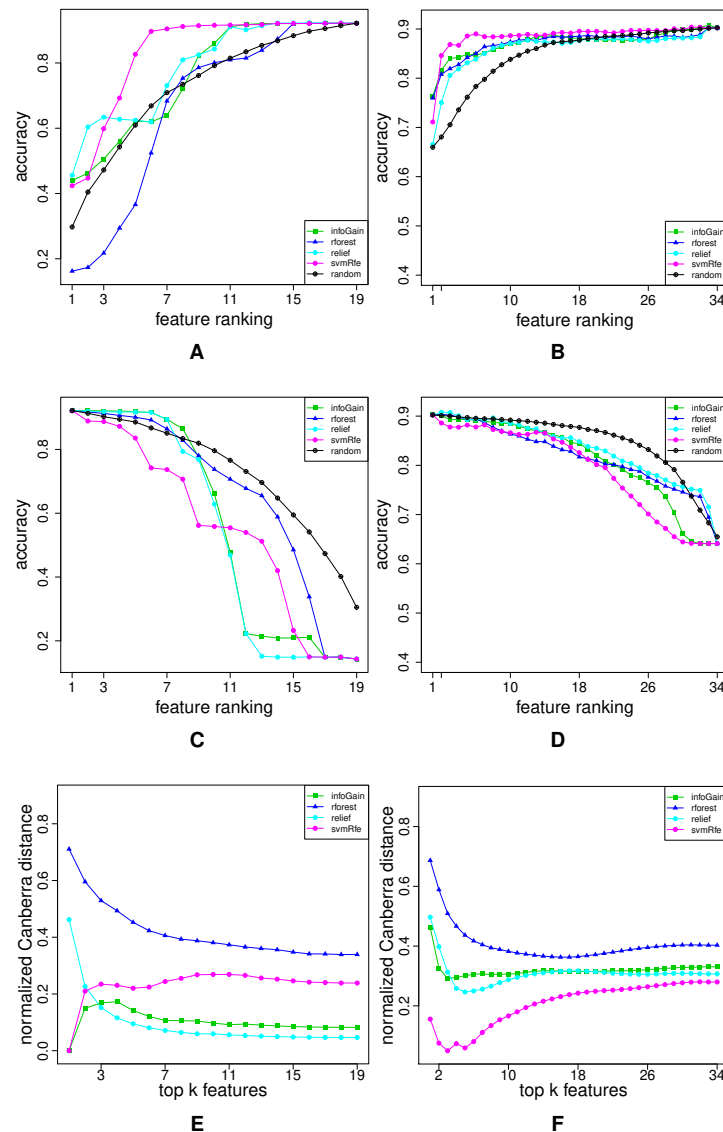




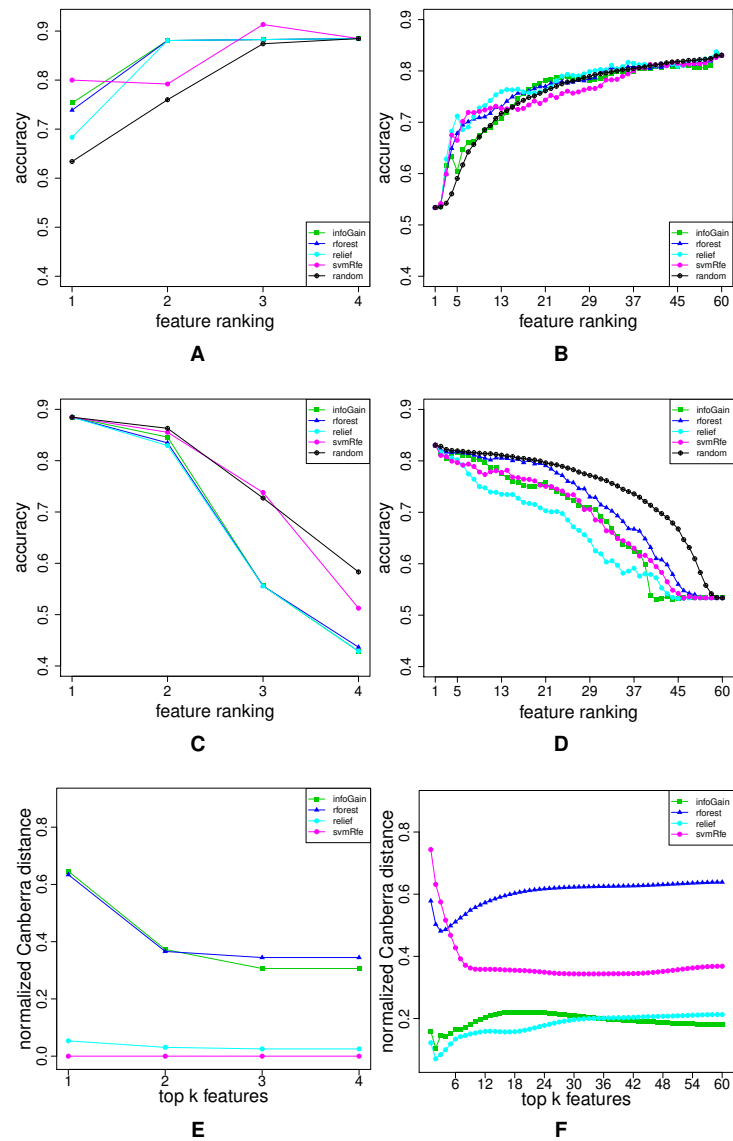
**Figure A7.** Ranking quality assessment for datasets `heart-c` (first column) and `heart-h` (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



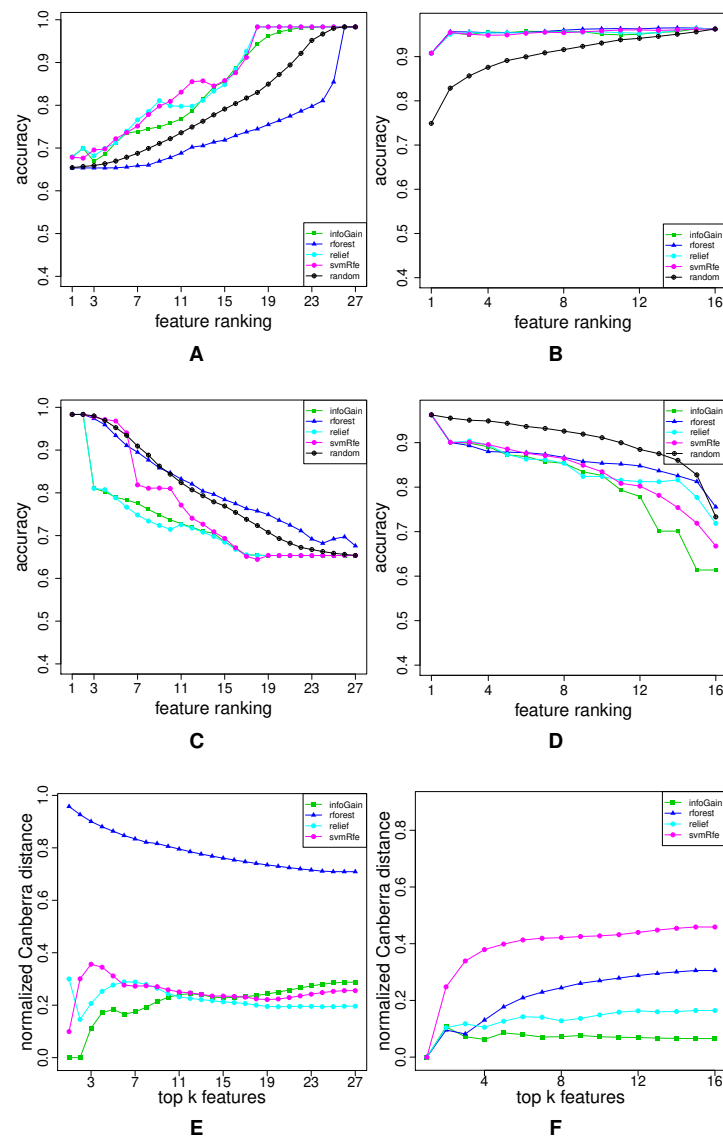
**Figure A8.** Ranking quality assessment for datasets `heart` (first column) and `hepatitis` (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



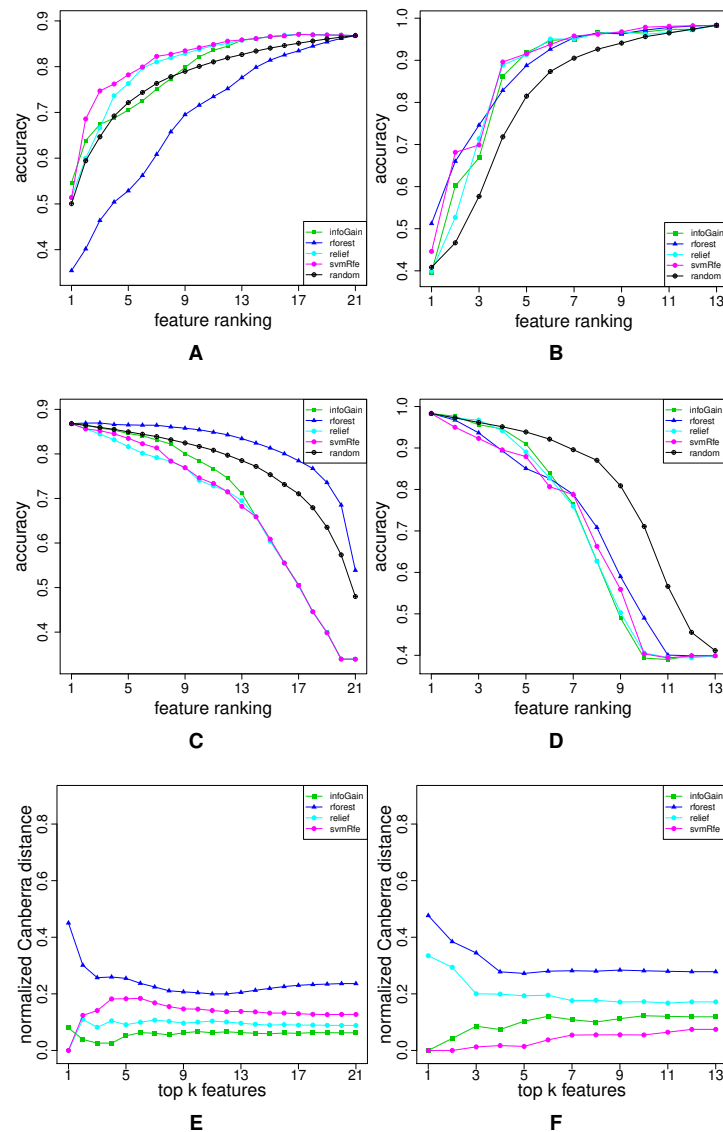
**Figure A9.** Ranking quality assessment for datasets *image* (first column) and *ionosphere* (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



**Figure A10.** Ranking quality assessment for datasets *iris* (first column) and *sonar* (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



**Figure A11.** Ranking quality assessment for datasets `tic-tac-toe` (first column) and `vote` (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.



**Figure A12.** Ranking quality assessment for datasets waveform (first column) and wine (second column) in terms of the FFA (first row) and RFA curves (second row), and rankings' stability estimates (third row). The FFA/RFA curves are obtained by using SVMs.