# An Empirical Comparison of Learning Algorithms for Nonparametric Scoring

Marine Depecker, Stéphan Clémençon, Nicolas Vayatis

**DRAFT manuscript No.**
(will be inserted by the editor)

# An Empirical Comparison of Learning Algorithms for Nonparametric Scoring

## The TREERANK Algorithm and Other Methods

**Stéphan Clémençon · Marine Depecker · Nicolas Vayatis**

the date of receipt and acceptance should be inserted later

Stéphan Clémençon · Marine Depecker
Télécom ParisTech
LTCI - UMR CNRS 5141
Paris, France
Tel.: +33-1-45817807
Fax.: +33-1-45817158
E-mail: stephan.clemencon@telecom-paristech.fr

Nicolas Vayatis
ENS Cachan & UniverSud
CMLA - UMR CNRS 8536 Cachan, France

# An Empirical Comparison of Learning Algorithms for Nonparametric Scoring

## The TREERANK Algorithm and Other Methods

**Abstract** The TREERANK algorithm was recently proposed in [1] and [2] as a scoring-based method based on recursive partitioning of the input space. This tree induction algorithm builds orderings by recursively optimizing the Receiver Operating Characteristic (ROC) curve through a one-step optimization procedure called LEAF-RANK. One of the aim of this paper is the in-depth analysis of the empirical performance of the variants of TREERANK/LEAFRANK method. Numerical experiments based on both artificial and real data sets are provided. Further experiments using resampling and randomization, in the spirit of bagging and random forests are developed ([3], [4]) and we show how they increase both stability and accuracy in bipartite ranking. Moreover, an empirical comparison with other efficient scoring algorithms such as RANK-BOOST and RANKSVM is presented on UCI benchmark data sets.

**Keywords** Scoring rules · Ranking trees · ROC curve · AUC Maximization · Resampling · Feature randomization

## 1 Introduction

In the last decade, learning to rank signals, web pages, individuals, organizations, has affected critical aspects of modern society. But behind the very concept of 'ranking', very diverse setups and numerous types of data may be considered. One of the subproblems in the ranking literature which has drawn a lot of attention is known as ranking with bipartite feedback or bipartite ranking [5]. In short, bipartite ranking can be cast as the problem of ordering instances $x$ of a possibly high dimensional space $\mathcal{X}$ based on a binary feedback information $y \in \{-1, +1\}$ attached to every instance $x$. The type of applications we have in mind are those related to scoring high dimensional observation vectors with binary labels, such as medical diagnosis or credit-risk screening. Several methods and algorithms can be applied to produce scoring rules for bipartite ranking: plain (or Kernel) LOGISTIC REGRESSION [6], [7], LOGITBOOST [8], RANKBOOST (see [5]), RANKSVM [9], RANKLS [10], RANKNET (see [11]), and many others. In the present paper, we will mainly consider a promising tree-based procedure called TREERANK recently developed and described in [1], [2]. We will explain how it differs from other algorithms in its very principles and explore its performance by providing an intensive experimental study, jointly with a fair comparison with other methods available in open source programming environments.

The guide for this study is the type of metrics used for performance assessment. In signal detection or medical diagnosis, the most complete description of the performance of a real-valued scoring rule is the ROC curve which plots the true positive rate as a function of the false positive rate as the detection threshold varies. However, as a function-like criterion, it only provides a partial order on the class of all possible scoring rules and it is quite difficult

to optimize in practice. The mostly used reference performance measures are summaries of the ROC curve and the most popular certainly is the Area Under the ROC Curve, or AUC. The drawback of this type of measure lies in the fact that it is a global measure and does not take into account the errors among the top of the list (highest scores). Various approaches for providing local measures of scoring performance related to the ROC curve have been proposed [12], [13], [14] or [15] and such metrics will be considered for assessing performance. There are numerous different criteria such as the Discounted Cumulative Gain, the Hit Ratio@x%, the Average Precision@n, and their variants (we refer to [16] for an extensive enumeration and discussion regarding all these criteria). Interestingly, the main feature of modern learning algorithms consists of performance maximization (or risk minimization, equivalently) and one expects that the closer the optimization risk functional is to the actual target risk functional, the better the performance evaluated on test data sets. This observation is the foundation for many machine learning algorithms extending the principle of convex risk minimization developed in the classification setup to the case of data conceived as pairwise comparisons [17]. The answer to the standard questions of regularized learning methods (choice of cost functions, penalty or dynamics - boosting or SVM-type) then determines specific ranking algorithms. For instance, RANKBOOST implements AUC maximization in the spirit of ADABOOST, RANKSVM proposes to fit an SVM algorithm with the hinge loss, RANKLS uses the least-squares as a cost function, and RANKNET follows the spirit of Neural Networks with an entropy-based cost function. The common feature of these algorithms is that they rely on pairwise comparisons (preferences) between labels and scoring functions. We emphasize the fact that the principle of performance maximization actually differs from the classical statistical methods used in scoring applications. Indeed, it appears that the mainstream scoring method of logistic regression is based on the conditional maximum likelihood maximization [6] and aims at modeling the unknown posterior probability

$$\eta(x) = \mathbb{P}\{Y = +1 \mid X = x\} .$$

This makes sense since the function $\eta$ actually defines the optimal scoring rule in the sense of the ROC curve [1]. Besides, statistical learning theory has established connections between convex risk minimization with classification - calibrated costs and logistic regression ([8]) and this idea has been further developed in kernel logistic regression algorithms ([7] and references therein). In these algorithms, the parametric estimation step of logistic regression is replaced by nonparametric estimation where an SVM-type convex optimization problem is solved. Eventually, an estimate of the function $\eta$ is provided by plugging-in the solution of the optimization step in a closed-form expression for $\eta$. Such methods seem to offer reasonable competitors to RANKBOOST. However, from the theoretical side, they present serious drawbacks. Indeed, in case that the statistical model used to represent the regression function is wrong, there are no guarantees about the stability of the estimated scoring rule. Such *plug-in* methods suffer from the curse of dimensionality and convex risk minimization methods are known to fail in the estimation of the regression function ([18],[19]). Hence, methods based on performance, such as RANKBOOST offers more guarantees than statistical estimation methods in leading to efficient ranking rules as their principle relies on the very optimization of one of the objective performance metrics for bipartite ranking. On the other hand scoring-based methods provide a richer output as they give a finer local description of the 'ranking'. Indeed, the limitation of RANKBOOST lies in the fact that AUC is a global criterion. Therefore there is no reason why methods like RANKBOOST should perform well on, say, the top 10% of instances in $\mathcal{X}$.

In the present paper, we propose to explore the empirical performance of TREERANK and a series of its variants (alternative splitting rule called LEAFRANK, and RANKING FORESTS, see [1], [2], [4]). The interesting feature of the TREERANK algorithm is that it takes the best of the two approaches previously described. Indeed, the TREERANK procedure uses decision trees to build a scoring rule which estimates the optimal ordering as induced by the regression function $\eta$, while maximizing the *Area Under the* ROC *Curve* (AUC) in a recursive manner. The one-step AUC maximization is implemented with a splitting rule called LEAFRANK which can be seen as an algorithm resolving a cost-sensitive

classification problem with data-dependent cost (see [2]). The reason why this method can potentially perform better than RankBoost is that TreeRank/LeafRank method optimizes the ROC curve in a stronger sense than the AUC (in sup-norm namely, see [1]). It also avoids the drawbacks of statistical estimation methods as it does not aim at the direct estimation of the regression function (only the ordering induced matters here, and not the actual values of the scoring rule).

The structure of the paper is the following. In Section 2, the setup and the notations are introduced. Section 3 presents the ranking algorithms under study in this paper, including TreeRank/LeafRank, Ranking Forests, and competitors such as RankBoost, RankLS, or RankSVM. The main results can be found in the numerical experiments described in Section 4. The conclusions are gathered in the discussion of Section 5.

## 2 The scoring problem in presence of classification data

We first introduce concepts from probabilistic modeling which describe the nature of the bipartite ranking problem. We shall refer here to the optimal elements for bipartite ranking and to the performance measures that can be retained to capture the optimal decision rules.

### 2.1 Setup

*Probabilistic model.* The probabilistic setup for bipartite ranking is the same as for binary classification. The data appear as a sample $\mathcal{D}_n = \{(X_i, Y_i) : 1 \leq i \leq n\}$, and we assume that the $(X_i, Y_i)$'s are *i.i.d.* realizations of a random pair $(X, Y)$, where X lies in a feature space $\mathcal{X}$ and $Y$ is a binary label in $\{-1, +1\}$. The joint probability distribution of the pair $(X, Y)$ can be described, for instance, by $P_X$ the marginal distribution of X and the regression function $\eta(x) = \mathbb{P}(Y = +1 \mid X = x)$, $\forall x \in \mathcal{X}$. We also use the notation $p = \mathbb{P}(Y = +1)$.

*Optimal scoring rules.* It is well-known that the target of binary classification is the level set $R^* = \{x \in \mathcal{X} : \eta(x) > 1/2\}$ ([20]). Indeed, predicting $y = +1$ on any $x \in R^*$ and $y = -1$ otherwise is the best possible decision in this

case. Now, if the problem is how to order the instances of $\mathcal{X}$ in such a way that positively labeled data appear on top of the list with high probability, then the optimal decision rules consist of real-valued scoring functions which imitate the ordering induced on $\mathcal{X}$ by the regression function $\eta$. In previous works ([1],[21]), we have showed that the following statements are true:

- *The class of optimal scoring functions are strictly increasing transforms of the regression function.*
  They are of the form: $s^* = T \circ \eta$, for some measurable function $T : \mathbb{R} \to \mathbb{R}$, whose restriction to the range of $\eta(X)$ is strictly increasing.
- *Optimal scoring rules can be obtained from the entire collection of level sets of the regression function $\eta$.*
  Indeed, they can be represented as: $\forall x \in \mathcal{X}$, $s^*(x) = c + \int_0^1 w(u) \cdot \mathbb{I}\{\eta(x) > u\} \, du$, where $c \in \mathbb{R}$, $w > 0$.
- *Bipartite ranking may be viewed as a continuum of classification problems.*
  This fact follows from the representation of $s^*$. Finding optimal scoring rules boils down to recovering the whole class of level sets: $\{\{\eta(x) > u\}, u \in [0, 1]\}$.
- *The* ROC *curve fully describes the performance of scoring functions in the bipartite ranking setup.*
  As a matter of fact, for a given distribution $P_X$, optimal scoring rules $s^*$ have an ROC curve above the ROC curve of any other scoring rule $s$. An easy bipartite ranking problem corresponds to the case where class-conditional distributions (distribution of $X$ given the class label $Y$) are far from each other. On the opposite, in the case when there is an important overlap between the two class-conditional distributions, the optimal ROC curve becomes close to the diagonal line.

### 2.2 Performance measures for scoring

Most of machine learning algorithms involve a functional optimization step. One of the main contributions of statistical learning theory has been the understanding of the connection between the optimization performed and the performance metrics for assessing the quality of

prediction ([17]). The choice of an objective criterion determines the optimal elements for a given problem and allows to have a baseline to compare prediction results for various methods. In the case of binary classification, considering that performance of a classifier $g : \mathcal{X} \to \{-1, +1\}$ is measured with the misclassification rate $L(g) = \mathbb{P}\big(Y \neq g(X)\big)$ implies that the optimal classifier is the indicator of the level set $R^*$. Now, if the criterion puts asymmetric weights on the two types of error

$$L_w(g) = 2w\mathbb{P}\big(Y = +1, g(X) = -1\big)$$
$$+ 2(1-w)\mathbb{P}\big(Y = -1, g(X) = +1\big) ,$$

then the optimal classifier is the indicator of a level set of the regression function whose level is equal to $w$. In addition, most of the efficient algorithmic approaches, such as boosting or Support Vector Machines, actually optimize surrogate (convex) risk functionals. Under rather weak assumptions, optimizing the surrogate criterion actually optimizes the misclassification rate and proving this fact was definitely one of the major achievements in the recent years (see [22]).

*The* ROC *curve and the* AUC. When it comes to scoring applications, the golden standard definitely is the ROC curve ([23]) as it represents the discrimination ability of a scoring rule at all levels. The ROC curve of a scoring rule $s$ plots the true positive rate $\beta_s(t) = \mathbb{P}\big(s(X) > t \mid Y = +1\big)$ against the false positive rate $\alpha_s(t) = \mathbb{P}\big(s(X) > t \mid Y = -1\big)$ as the threshold $t$ varies from $-\infty$ to $+\infty$. However, the ROC curve is a complex object to optimize as it represents a function-like criterion. It mainly serves as a visual display of performance used *a posteriori* (i.e. after training the scoring rule). A simpler object is the functional defined by the Area Under the ROC Curve, known as the AUC which can be interpreted as the probability of concordant pairs:

$$\text{AUC(s)} = \mathbb{P}\big(s(X) > s(X') \mid Y = 1, Y' = -1\big)$$
$$+ \frac{1}{2}\mathbb{P}\big(s(X) = s(X') \mid Y = 1, Y' = -1\big),$$

where $(X, Y)$ and $(X', Y')$ are i.i.d. pairs of observations.

*Remark 1* In the case of preference-based ranking, then individual labels are not available and

there exists no ROC curve for a preference function. However, it is possible to extend the definition of AUC to this setup. Indeed, denote preference data by triples of the form $(X, X', Z)$ where $X, X' \in \mathcal{X}$, and $Z \in \{-1, 0, +1\}$. Then a preference function $\pi : \mathcal{X} \times \mathcal{X} \to \{-1, 0, +1\}$ has an AUC defined by

$$\text{AUC}(\pi) = \mathbb{P}\big(\pi(X, X') > 0 \mid Z = +1\big)$$
$$+ \frac{1}{2}\mathbb{P}\big(\pi(X, X') > 0 \mid Z = 0\big) .$$

Note that a scoring rule $s$ leads to the obvious choice for a preference function : $\pi(x, x') = s(x) - s(x')$ (see [24] for details and optimal elements for preference data). The converse is less straightforward (see for example [25]).

An important limitation of the AUC is that it only measures the quality of a scoring / ordering in a global manner. Two scoring rules may present the same AUC and simultaneously have a very different behavior. Indeed, the relative weight of a discordant pair $(x, x')$ is only a function of the difference $s(x) - s(x')$ of the scores between the two instances and not of their position over the range of $s$. This is a crucial remark, especially as in most applications involving ranking problems people often focus on the instances with highest scores (corresponding to the part of the ROC curve close to the origin). As highlighted in [13], page ranking, credit risk screening or medical diagnosis are good examples of fields in which priority is given to observations among the top scores. Indeed, in credit risk screening for instance, the main objective will be to identify potential risks among the best candidates for a credit. In a similar way, main objective in medical diagnosis will be to ensure correct diagnosis among the more risky patients of a cohort. There exist several criteria which emphasize the highest scores : the Average Precision (AP), the Predicted-Rank-of-Top (PROT), the Coverage, the Discounted Cumulative Gain (DCG), the $p$-norm push ([15]). Most of these criteria can be related to linear rank statistics as shown in [13] and [14], but they are not necessarily related to specific algorithms maximizing them.

*Local* AUC. In [13], a truncation of the AUC called the local AUC has been introduced which examines performance in the sense of the AUC

but only at the top of the ranked list. The truncation is shown to be consistent with the bipartite ranking problem as the regression function $\eta$ or any strictly increasing transform of it is still optimal. Fix $u \in (0,1)$ the rate of best instances and let $t^* = t(s,u)$ be the corresponding threshold: $u = \mathbb{P}\big(s(X) > t^*\big)$. We take the notations: $\alpha^* = \alpha_s(t^*)$ and $\beta^* = \beta_s(t^*)$, which are determined as the coordinates of the intersection of the *control line* $d_u : u = p\beta + (1-p)\alpha$ in the ROC space with coordinates $(\alpha, \beta)$ with the ROC curve (see Figure 1). Then, the Local AUC is defined by:

$$\mathrm{LocAUC(s,u)} = \beta^*(1 - \alpha^*)$$
$$+ \int_{t=0}^{\alpha^*} \beta_s \circ \alpha_s^{-1}(t) \ dt,$$

where $\alpha_s^{-1}(t) = \inf\{u \in (0,1) : \alpha_s(u) \geq t\}$ for all $t \in (0,1)$. Instead of putting more weight on top-ranked instances, the local AUC only focuses on the top of the list with a truncation.
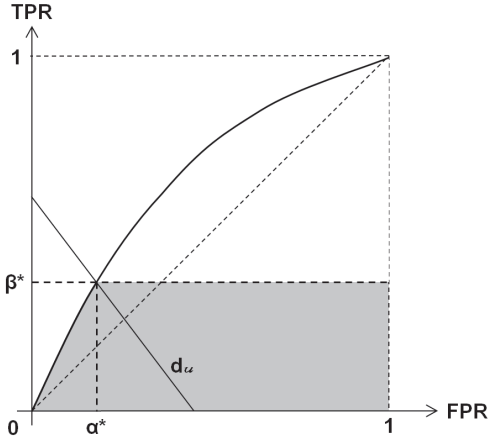


**Fig. 1** The ROC curve, the control line and the LocAUC criterion.

*Other measures.* The tree-based approach considered in this paper focuses on the ROC curve and the AUC. In the experimental section, we will also consider other criteria for completeness. We shall provide the Average Precision as well as the Hit Rate.

## 3 TreeRank and other scoring methods

Our focus in the paper lies in the scoring-based method proposed in [1] and [2], using recursive partitioning in the spirit of decision trees. In the present section, we introduce the main ideas and the algorithms in this approach, as well as some of its competitors.

### 3.1 Presentation of TreeRank and its variants

*Ranking trees.* Scoring rules output by TreeRank/LeafRank can be described as oriented binary trees, called ranking trees. Ranking trees are defined by two elements: (i) a partition of the input space $\mathcal{X}$ into cells, and (ii) a permutation on the cells which indicates the ordering induced by the rule. Denote by $\mathcal{T}_D$ the binary tree of depth $D \geq 0$ which represents a recursive partition of the input space $\mathcal{X}$. Let us equip the tree with a left-to-right orientation. We then obtain a collection of scoring rules stored in a tree-like structure (see Figure 2). The root of the tree is $C_{0,0} = \mathcal{X}$, and each internal node $C_{d,k}$, with $0 \leq d < D$ and $0 \leq k < 2^d$, corresponds to a cell in $\mathcal{X}$ and splits into two non-empty cells $C_{d+1,2k}$ (left sibling) and $C_{d+1,2k+1} = C_{d,k} \setminus C_{d+1,2k}$ (right sibling). Any subtree $\mathcal{T} \subset \mathcal{T}_D$ can then be used to define a pre-order on $\mathcal{X}$: all elements in the same terminal cell will be considered as ties, and the ordering is induced by the left-to-right orientation of the subtree. The scoring rule is then defined as:

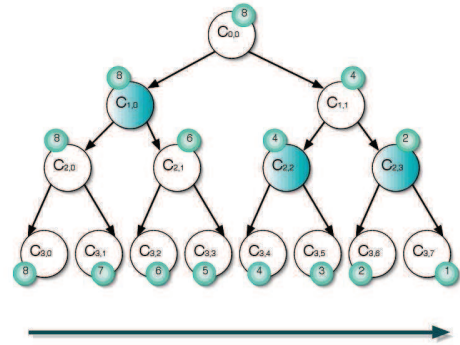$$s_{\mathcal{T}}(x) = \sum_{C_{d,k}:\text{terminal cell}} 2^D(1 - k/2^d) \cdot \mathbb{I}\{x \in C_{d,k}\}.$$



**Fig. 2** Ranking tree.

*Recursive partitioning for ranking.* The TreeRank algorithm has been recently introduced

in [1] and its theoretical properties have been discussed there in detail. We now recall the principle of this method. The algorithm follows the principle of CART proposed in [26], excepts that it builds an oriented binary tree and uses a splitting rule which implements AUC maximization at the node level (optimization step). Consider $(X_1, Y_1), \ldots, (X_n, Y_n)$, a sample of classification data. One of the findings of the analysis is that the optimization step corresponds to solving a cost-sensitive classification problem. At a given node, the cell $C$ is split into two parts $C_+$ and $C_- = C \setminus C_+$ so that the weighted empirical error

$$\widehat{L}_{C,\omega}(\Gamma) = \frac{2(1-\omega)}{n} \sum_{i=1}^{n} \mathbb{I}\{X_i \in C \setminus \Gamma\} \mathbb{I}\{Y_i = 1\}$$

$$+ \frac{2\omega}{n} \sum_{i=1}^{n} \mathbb{I}\{X_i \in \Gamma \cap C\} \mathbb{I}\{Y_i = -1\} \ .$$

is minimized at $\Gamma = C_+$, for an adaptively chosen value of the weight $\omega$. Then the TREERANK algorithm for the recursive optimization of the ROC curve can be formulated as follows. We point out that the class of candidate sets $\Gamma$ is determined by the choice of a splitting rule.

---

THE TREERANK ALGORITHM

- **Initialization.** Start with $C_{0,0} = \mathcal{X}$.

- **Itérations.** For $d = 0, \ldots, D - 1$ and $k = 0, \ldots, 2^d - 1$ ,

  1. Set $n_{d,k} = \sum_{i=1}^{n} \mathbb{I}\{X_i \in C_{d,k}\}$. Compute the rate of false positives in cell $C_{d,k}$:

  $$\alpha(C_{d,k}) = \frac{1}{n_{d,k}} \cdot$$
  $$\sum_{i=1}^{n} \mathbb{I}\{X_i \in C_{d,k}\} \cdot \mathbb{I}\{Y_i = +1\}.$$

  2. Minimize $\widehat{L}_{C_{d,k},\omega}(\Gamma)$ with

  $$\omega = \omega_{d,k} = \alpha(\mathcal{C}_{d,k}) \ .$$

  3. Set $C_{d+1,2k} = C_+$ the minimizer, and $C_{d+1,2k+1} = C_{d,k} \setminus C_{d+1,2k}$.

- **Output.** Oriented binary tree

  $$\mathcal{T}_D = \{C_{d,k} : 0 \le d \le D, \ 0 \le k < 2^d\}.$$

---

*Splitting rules.* The choice of a splitting rule relies on the the trade-off between efficiency and computability. In plain decision trees, simple rules are retained such as orthogonal splits. Here the task being more involved, more complex splitting rules may be considered. In principle, any classification algorithm with an easy tuning of asymmetric costs can do the job. In ([2]), a specific splitting rule is proposed, called LEAFRANK which uses the TREERANK algorithm at the node level recursively to build a partition and then splits the cells of this partition into two categories: one collection of cells going to the left sibling node, and the rest going to the right sibling node. The classification of cells relies on a result which states that, for a fixed partition, the optimal permutation in the sense of the AUC is the one corresponding to a monotone ordering of the ratio of true positive rate over false positive rate per cell. The advantage of LEAFRANK is that it uses orthogonal splits and hence preserves the interpretability of the scoring rule. It is possible indeed to recover measures of importance of the variables and of dependence between them (see [2]). An alternative to LEAFRANK is to use a cost-sensitive version of a SVM classifier (see [27]), which increases the complexity of the scoring rule but at the same time, interpretability of features is lost.

*Pruning.* Decision trees are known to suffer from overfitting as you can end up with one observation per cell. Beyond using arbitrary stopping rules, more can be done with pruning. Pruning consists in removing subtrees from the complete tree grown with a given splitting rule. This can be done in several ways and we consider here the simple option of taking a linear penalty with a constant calibrated through cross validation. The penalized criterion is the following:

$$\widehat{\text{AUC}}(s_\tau) + \lambda \cdot \#\mathcal{P}(\tau) \ ,$$

where $\tau$ is a pruned tree, $s_\tau$ is the corresponding scoring function, $\widehat{\text{AUC}}(s)$ is the empirical counterpart of the AUC, $\#\mathcal{P}(\tau)$ is the cardinal of the partition induced by $\tau$ on the input space $\mathcal{X}$, and $\lambda$ is a constant to be estimated through cross validation.

*Resampling and randomization.* Aggregation procedures in machine learning are known to improve both performance and robustness of de-

cision rules. This principle underlies numerous algorithms derived from bagging and boosting methods. Using the idea of resampling and aggregation, a collection of scoring rules can be built from many bootstrap samples in the spirit of bagging methods. By adding the ingredient of randomization an analogue to random forests called Ranking Forests can be derived (a presentation is given in [4]). For completeness, its description is provided in the box below.

*The force of* TREERANK. The strengths of the TREERANK algorithm are numerous. First, it possesses the main features of tree-based decision rules from the algorithmic point of view. It is also straightforward to read the contribution of latent variables when using decision trees as splitting rules at each node. One ranking tree by itself also suffers the same drawback as CART, and that is instability. Variants based on aggregation of many ranking trees, such as bagging and ranking forests, lead to more robust methods and improve performance significantly as shown in the experimental section. Another important feature of TREERANK lies in its theoretical soundness. Indeed, it can be shown (see [1]) that the scoring rule output by TREERANK produces near-optimal ROC curves in the supnorm sense.

### 3.2 Other learning methods for ranking and scoring classification data

The optimization principles underlying several algorithms used in bipartite ranking are of two types: those which focus on the optimization of the target performance measure and those which aim at recovering the target regression function $\eta(x) = \mathbb{P}\{Y = +1 \mid X = x\}, \forall x \in \mathcal{X}$. In the first category belong methods with an optimization based on a pairwise criterion. This is indeed the case with AUC maximization methods such as RANKBOOST, RANKSVM, RANKLS. Other algorithms like RANKNET optimize a nonconvex criterion but they also use pairwise cost functions. Among the methods which focus on the optimal function to recover, we consider classification methods which output a real-valued scoring function such as LOGISTIC REGRESSION, ADABOOST and LOGITBOOST, or KERNEL LOGISTIC REGRESSION among others.

RANKBOOST ([5]). This method mimics ADABOOST in the pairwise setup. It boils down

to the optimization of the following criterion

$$\frac{2}{n(n-1)} \sum_{i<j} \exp\{(Y_i - Y_j)(f(X_i) - f(X_j))\} \, ,$$

$f$ being a linear combination of weak ranking rules. Weak ranking rules involved are decision stumps.

The functional above is a proxy to the empirical AUC.

RANKSVM and RANKLS (resp. [9] and [10]). These two methods are export the SVM methodology to set a pairwise criterion of the form

$$\frac{2}{n(n-1)} \sum_{i<j} d\big(Y_i - Y_j, f(X_i) - f(X_j)\big) + \lambda \|f\|_K,$$

where $f(x) = \sum_{i=1}^{n} \alpha_i K(x, X_i)$, $\lambda$ is a smoothing parameter and $\|f\|_K$ is the RKHS norm of $f$ related to the kernel $K$.

In the case of RANKSVM the loss function $d$ is the hinge loss $d(u, v) = (1 - uv)_+$, while for RANKLS, it corresponds to the least-squares criterion $d(u, v) = (u - v)^2$.

RANKNET ([11]). Here the cost function used is based on cross entropy between the target value of the probability of a pair and the estimated probability

$$d\big(Y_i - Y_j, f(X_i) - f(X_j)\big) = -\Delta_{i,j} \log \pi_{i,j}$$
$$-(1 - \Delta_{i,j}) \log(1 - \pi_{i,j}),$$

where

$$\Delta_{i,j} = \frac{1 + (Y_i - Y_j)/2}{2} \in \{0, 1/2, 1\} \, , \text{ and}$$
$$\pi_{i,j} = \pi(f(X_i) - f(X_j)) = \frac{e^{f(X_i) - f(X_j)}}{1 + e^{f(X_i) - f(X_j)}} \, .$$

and the candidate functions $f$ have the same structure as neural nets.

Note that the cost function can also be written as:

$$d\big(Y_i - Y_j, f(X_i) - f(X_j)\big) \;=\;$$
$$\log\left(1 + e^{f(X_i) - f(X_j)}\right) - \Delta_{i,j} \left(f(X_i) - f(X_j)\right).$$

LOGITBOOST ([8]) and KERNEL LOGISTIC REGRESSION ([7]).

These two methods do not involve pairs of observations and therefore belong to classification methods, except that they output a real-valued

8

---

<div style="border:1px solid black">

THE RANKING FOREST ALGORITHM

1. **Parameters.**

   - $B$ number of bootstrap replicates,

   - $n^*$ bootstrap sample size,

   - tuning parameters of TREERANK - depth $D$ and presence/absence of pruning,

   - $\mathcal{FR}$ feature randomization strategy,

   - $d$ pseudo-metric over orderings.

2. **Bootstrap profile makeup.**

   (a) (RESAMPLING STEP.) Build $B$ independent bootstrap samples $\mathcal{D}_1^*$, ..., $\mathcal{D}_B^*$, by drawing with replacement $n^* \cdot B$ pairs among the original training sample $\mathcal{D}_n$.

   (b) (RANDOMIZED TREERANK.) For $b = 1$, ..., $B$, run TREERANK combined with the feature randomization method $\mathcal{FR}$ based on the sample $\mathcal{D}_b^*$, yielding the ranking tree $\mathcal{T}_b^*$, related to the partition $\mathcal{P}_b^*$ of the space $\mathcal{X}$.

   The aggregation results from a median ranking rule over the largest subpartition $\mathcal{P}^*$ extracted from the collection of partitions implemented by the ranking trees. Indeed, every ranking tree $\mathcal{T}$ induces an ordering $\preceq_\mathcal{T}$ over the cells of $\mathcal{P}^*$ and the median ranking tree induces an ordering defined as the solution of a distance minimization problem.

   $$\preceq^* = \arg\min_{\preceq \in \mathbf{R}(\mathcal{P}^*)} \sum_\mathcal{T} d(\preceq, \preceq_\mathcal{T}) \ ,$$

   where $\mathbf{R}(\mathcal{P}^*)$ denotes all possible orderings of cells in the subpartition $\mathcal{P}^*$, and $d$ defines a distance between orderings, such as the Kendall tau (count of concordant pairs when comparing the two orderings), Spearman rule ($L_1$ distance between the orderings seen as vectors) or Spearman correlation coefficient ($L_2$ distance between the orderings), for instance.

3. **Aggregation.** Compute the largest subpartition partition $\mathcal{P}^* = \bigcap_{b=1}^B \mathcal{P}_b^*$. Let $\preceq_b^*$ be the ranking of $\mathcal{P}^*$'s cells induced by $\mathcal{T}_b^*$, $b = 1$, ..., $B$. Compute a median ranking $\preceq^*$ related to the bootstrap profile $\Pi^* = \{\preceq_b^*: 1 \le b \le B\}$ with respect to the metric $d$ on $\mathbf{R}(\mathcal{P}^*)$:

   $$\preceq^* = \arg\min_{\preceq \in \mathbf{R}(\mathcal{P}^*)} \sum_{b=1}^B d(\preceq, \preceq_b^*) \ ,$$

   as well as the scoring rule $s_{\preceq^*, \mathcal{P}^*}(\mathrm{x})$.

</div>

function which, in theory, is related to the target regression function $\eta$. The criterion used in both methods is the following:

$$\frac{1}{n} \sum_{i=1}^n \log\left(1 + \exp\{-Y_i \cdot f(X_i)\}\right) \ ,$$

or (if $Y_i \in \{0, 1\}$)

$$\frac{1}{n} \sum_{i=1}^n \left(-Y_i \cdot f(X_i) + \log\left(1 + \exp\{f(X_i)\}\right)\right) \ .$$

An iterative boosting-type procedure is used for the optimization in LOGITBOOST, while KER-

NEL LOGISTIC REGRESSION relies on convex risk minimization with a penalty proportional to the RKHS norm of $f$.

P-NORM PUSH ([15]). The criterion used in this method is very similar to RANKBOOST but the idea is to put more weight on the largest discrepancies. The algorithm implements the optimization of the following criterion:

$$\frac{2}{n(n-1)} \cdot \sum_{i=1}^n$$

$$\left(\sum_{j=i+1}^n \exp\left\{(Y_i - Y_j)(f(X_i) - f(X_j))\right\}\right)^p$$

with $p \geq 1$, and $f$ being a linear combination of weak ranking rules as in RANKBOOST. We point out that in the case where $p = 1$, the P-NORM PUSH is the same algorithm as RANK-BOOST, larger values of $p$ put more emphasis on the largest scores.

In the present paper, we only consider those methods for which there is a publicly available implementation: KLR, RANKBOOST, RANKLS, RANKSVM, ADABOOST and also P-NORM PUSH.

## 4 Numerical experiments

### 4.1 Data sets

The methods presented in the paper have been tested on several data sets, both real (10) and artificial (3).

- Artificial data sets are gaussian mixtures, in the space $\mathbb{R}^{20}$, called *GaussEasy20d*, *Gauss-Med20d* and *GaussHard20d*. They have an increasing complexity as the distance between the centers decreases. The covariance matrices for each class conditional distributions can be different but both are diagonal.

- Real data sets are ten benchmark data from the UCI repository and available at

  <http://archive.ics.uci.edu/ml/>

  with their characteristics given in Table 1.

### 4.2 Description of algorithms

In the series of experiments presented in this section, several variants of TREERANK have been tested and compared among each other and to their main competitors. More precisely, we first compare twelve variants of TREERANK on the three simulated toy examples. We then retain eight of them for an insightful comparison on the benchmark datasets from the UCI repository. Eventually, we compare two well-performing versions of the TREERANK algorithm with the six methods previously cited : RANKBOOST, RAN-KLS, RANKSVM, KLR, ADABOOST and P-NORM PUSH.

Before specifying the TREERANK versions considered in these experiments, we first state some general observations and notations regarding the parametrization of these algorithms.

As explained in 3.1, any classification method can be used as *splitting rule* for growing a ranking tree, the choice of it mainly depending on the expectations of the user regarding the ranking rule (interpretability, computational cost, flexibility, etc.). In these experiments we have been using both the LEAFRANK splitting rule previously described and SVM classifiers. In the latter case, the kernel involved is either linear or gaussian, the one leading to the best performances being retained.

Besides, TREERANK variants based on re-sampling involve an aggregation procedure, leading to the *consensus* ranking (see 3.1). Here, we compute the median ranking rule using the pseudo-distance between orderings induced by the Spearman correlation coefficient.

Eventually, two additional parameters, $d_1$ and $d_2$, have to be considered when running the RANKING FOREST algorithm. Both parameters are related to feature randomization, affecting respectively the master ranking tree output by TREERANK and the subtrees output by the LEAFRANK procedure. Precisely, $d_1$ denotes the number of predictors randomly chosen (among the total number of dimensions) to optimize the local AUC at each node of the master ranking tree. At each step of the TREERANK procedure, the best split will be chosen among those $d_1$ predictors. Similarly, $d_2$ denotes the number of predictors randomly chosen at each node of the classification subtrees when the LEAFRANK procedure is used as splitting rule. Again, at each step of the latter, the best split will be chosen among those $d_2$ predictors. If we denote by $d$ the total number of explanatory variables (which is also the dimension of the observation vectors $X_i$), the following inequality stands true: $d_2 \leq d_1 \leq d$.

We now provide the specifications for each variant of the TREERANK algorithm used in the numerical experiments described below.

- TRK$_{\text{CART}}$ - single ranking tree with a splitting rule based on the LEAFRANK procedure,

- TRK$_{\text{SVM}}$ - single ranking tree with a splitting rule based on a SVM classifier,

| Name | Sample de Size | Number of variables | Rate of positives | Categorical variables | Description |
|---|---|---|---|---|---|
| *GaussEasy20d* | 2000 | 20 | 0.5 | no | gaussian mixture |
| *GaussMed20d* | 2000 | 20 | 0.5 | no | gaussian mixture |
| *GaussHard20d* | 2000 | 20 | 0.5 | no | gaussian mixture |
| *Australian Credit* | 690 | 14 | 0.44 | yes | Credit risk data |
| *German Credit* | 1000 | 20 | 0.7 | yes | Credit risk data |
| *Japanese Credit* | 690 | 15 | 0.45 | yes | Credit risk data |
| *Breast Cancer Diagnosis* | 569 | 30 | 0.37 | yes | Detection de malignant tumors |
| *Breast Cancer Original* | 683 | 9 | 0.35 | yes | Detection of malignant tumors |
| *Heart Disease* | 270 | 13 | 0.44 | yes | Detection of heart illness |
| *Hepatitis* | 112 | 18 | 0.17 | yes | Detection of hepatitis |
| *Autos MPG* | 392 | 7 | 0.53 | yes | Levels of fuel consumption |
| *Congressional Vote* | 232 | 16 | 0.53 | yes | American Congress voting results |
| *Ionosphere* | 351 | 34 | 0.64 | yes | Detection of structures |

**Table 1** Description of synthetic and UCI benchmark data.

- $\text{Bagg}_{\text{CART}}$ - bagged version based on the aggregation of $B = 50$ *resampled* ranking trees, which growth is based on the LEAF-RANK procedure,

- $\text{Bagg}_{\text{SVM}}$ - bagged version based on the aggregation of $B = 50$ *resampled* ranking trees, which growth is based on a SVM classifier,

- $\text{RF1}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor d/4 \rfloor$), which growth is based on the LEAFRANK procedure,

- $\text{RF1}_{\text{SVM}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor d/4 \rfloor$), which growth is based on a SVM classifier,

- $\text{RF2}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor 1 \rfloor$), which growth is based on the LEAFRANK procedure,

- $\text{RF2}_{\text{SVM}}$ - RANKING FORESTversion based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor 1 \rfloor$), which growth is based on a SVM classifier,

- $\text{RF3}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_2 = \lfloor d/4 \rfloor$), which growth is based on the LEAFRANK procedure,

- $\text{RF4}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_2 = \lfloor 1 \rfloor$), which growth is based on the LEAFRANK procedure,

- $\text{RF5}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor d/2 \rfloor$ and $d_2 = \lfloor d/4 \rfloor$), which growth is based on the LEAFRANK procedure,

- $\text{RF6}_{\text{CART}}$ - RANKING FOREST version based on the aggregation of $B = 50$ *resampled and randomized* ranking trees (with $d_1 = \lfloor d/4 \rfloor$

and $d_2 = \lfloor 1 \rfloor$), which growth is based on the LeafRank procedure.

### 4.3 Performance and stability metrics

*Performance.* In the following experiments we measured the algorithms performances mainly in terms of AUC and local AUC, but also with other criteria. When running the ranking algorithms on real datasets, these quantities are estimated using $V$-fold cross validation. We recall the principle of for a generic performance measure $\Pi$: from the training set $\mathcal{D}_n$, consider $V = 10$ blocks $\{\mathcal{D}^{(1)}, \cdots, \mathcal{D}^{(V)}\}$ of equal size which are randomly chosen and the expected performance is estimated as follows:

$$\overline{\Pi} = \frac{1}{V} \sum_{v=1}^{V} \widehat{\Pi}^{(v)},$$

where, for each $v \in \{1, \cdots, V\}$, the quantity $\widehat{\Pi}^{(v)}$ denotes the empirical estimate of $\Pi$ for the scoring rule $\widehat{s}^{(-v)}$ trained on the data set $\mathcal{D}^{(-v)} = \mathcal{D}_n \setminus \mathcal{D}^{(v)}$ and computed over $\mathcal{D}^{(v)}$. We apply the cross validation to the following criteria:

- AUC - Area Under an ROC curve,
- AUC@$u$% - Local AUC with values $u = 5, 10, 20\%$ for the proportion $u$ of best instances ([13]),
- AP - Average Precision,
- HR@$x$% - Hit Ratio with proportions $x = 10\%, 20\%$.

Additionally, in the case of the AUC and local AUC, we also provide the standard type error $\widehat{\sigma}^2$ of the collection of $V$-fold empirical estimates.

*Stability.* Stability is another interesting indicator to consider when comparing several algorithms. We propose two different and complementary ways to evaluate it. As $V$-fold cross validation is used on every dataset to provide *mean* values of the AUC and Local AUC, we provide for each algorithm an estimate of the standard deviation of the considered performance metrics. Moreover, on the simulated toy examples, we also consider an *instability* measure given by

$$\mathbf{Instab}(\mathbf{S}) = \mathbb{E}\left[d(\preceq_b, \preceq_{b'})\right],$$

where $d$ is a distance between orderings based on Spearman correlation coefficient and the orderings $\preceq_b$ and $\preceq_{b'}$ are respectively induced on the largest subpartition of the feature space $\mathcal{X}$ by two *resampled* ranking trees, meaning that the scoring rules $\preceq_b$ and $\preceq_{b'}$ are obtained by running the TreeRank heuristic on two *i.i.d.* learning datasets, resampled from the original one $\mathcal{D}_n$. Based on a resampling procedure, this indicator can be easily estimated through the following formula:

$$\widehat{\mathbf{Instab}} = \frac{2}{B(B-1)} \sum_{1 \leq b < b' \leq B} d(\preceq_b, \preceq_{b'}),$$

where $B$ is the total number of bootstrap samples drawn from $\mathcal{D}_n$.

### 4.4 Results

*Computational cost.* Before going into the detail of the experimental results, we give a first insight into the computational cost of the TreeRank variants. It is well known that tree-based algorithms can yield high computational costs, although in most cases, optimal trees are grown through bottom-up greedy procedures, avoiding a costly exhaustive search among the collection of optimal trees.

In our experiments, the computational cost mainly comes from the cross validation based pruning stage within both TreeRank and LeafRank procedures (see 3.1). Indeed, considering the *GaussEasy20d* dataset (2000 observations for 20 relevant dimensions), growing a ranking tree through the TRK$_{\text{CART}}$ variant takes less than 2 seconds when the sizes are settled as input parameters. Yet, it takes about 20 seconds when the size of both the subtrees and the master ranking tree are automatically selected through 5-fold cross validation. Thus, cross validation based pruning will yield heavy computational costs when running the Ranking Forests procedure, which will grow and aggregate $B = 50$ ranking trees. Indeed, growing a single ranking tree through Bagg$_{\text{CART}}$ or any other Ranking Forest variant on the *GaussEasy20d* dataset takes about 15 minutes with a 5-fold cross validation based pruning. This computational time depends on the randomization level: the more randomized the algorithm, the less costly. Yet, it can be reduced to 2 minutes when the

sizes of the grown trees are settled *a priori*.

Besides, making use of SVM based splitting rules yields heavier computational costs. Thus, the Bagg$_{\text{SVM}}$ procedure and its randomized counterparts need about 20 minutes to grow an automatically pruned (5-fold cross validation) aggregated ranking tree, and about 4 minutes when the sizes of the trees are settled in advance. Eventually, the computational costs observed on most of the benchmark datasets are lower to those detailed above (or comparable in some cases), due to the lower dimension and the sparsity of the problems.

*Comparison between the variants of* TREERANK. The main objective of this series of experiments is to compare the behavior of several versions of the TREERANK algorithm. We first run twelve heuristics on three toy examples based on gaussian mixtures: *GaussEasy20d, GaussMed20d* and *GaussHard20d*. In order to provide mean values for the performance indicators, we run all heuristics on $N = 30$ *i.i.d.* learning resamples of $\mathcal{D}_n$, with $n = 1000$ copies of the random pair $(X, Y) \in \mathcal{X} \times \{-1, +1\}$, and test the output scoring rules on a single test sample of $n = 1000$ observations. Given the computational costs described previously, for this series of experiments trees are not pruned automatically. For sake of computational time, the sizes of the grown trees (and subtrees) has been defined *a priori*: the maximum depth of the master ranking trees has been settled to $d_1 = 8$ and the maximum depth of the sub-ranking trees output by the LEAF-RANK procedure to $d_2 = 6$. Figure 3 gives a global insight into the performances of TREE-RANK variants on these toy examples. The average rank of each variant over the three toy examples is displayed for the performance indicators described in previous Subsection. Moreover, visual displays for each datasets are provided on Figure 4. The histograms representing both performance (AUC) and standard type error ($\hat{\sigma}^2$). We also report the Local AUC measure corresponding to a proportion $u = 20\%$ of best instances. Eventually, the quantitative results can be found in Tables 2, 3 and 4.

The overview given by Figure 3 shows that RF1$_{\text{SVM}}$ clearly outperforms all other procedures, except in terms of instability. We also notice that resampling and randomization yield



**Fig. 3** Average ranks of TREERANK variants on the three toy example for several performance indicators.

better results, both in terms of learning performance and robustness, while the level of randomization remains *reasonable*. Indeed, highly randomized procedures RF4$_{\text{CART}}$, RF5$_{\text{CART}}$ and RF6$_{\text{CART}}$ yield the worst results.

The graphs of Figure 4 allow to confirm our first observations and to go more into details. Thus, we notice that, when considering single ranking trees, the best performance is achieved by scoring rules using SVM classifiers as splitting rule, due to a superior representational capacity. As expected, performance is increased when aggregating many ranking trees. It is also interesting to notice that the choice of the splitting rule seems to have a weaker impact on the performances of those heuristics. Another interesting observation is related to the effect of feature randomization: it seems that randomization does not affect significantly the heuristics based on an SVM-type splitting rule. Moreover, as already mentioned, the more important the level of randomization, the worst the AUC performance (and the less robust the estimation), but we will see in the next experiment that this is not necessarily the case when considering real data sets. Actually, this phenomenon is not really surprising and can be explained in two different ways. As described in 3.1, RANKING FOREST is a combination of *resampling, randomization* and *aggregation*. Although it is now well known in Machine Learning that *aggregating bootstrap* versions of a classifier improves both its performance and robustness, it is also quite clear that picking up *randomly* a subset of predictors to learn a model introduces some instability. The underlying principle of RANKING

a. Sample GaussEasy20D



b. Sample GaussMed20D



c. Sample GaussHard20D

**Fig. 4** Comparison of TREERANK variants on simulated datasets.

FORESTS is precisely that randomization could improve the learning capacity of a classifier by focusing on subsets of most relevant predictors, the instability introduced in the process being compensated by the aggregation of several *re-sampled* classifiers. A first explanation to the results displayed on Figure 4 is that aggregating $B = 50$ ranking trees might not be sufficient to compensate the instability introduced by the different levels of randomization considered in the experiment. Moreover, it seems reasonable to state that randomization will perform better when the optimal model is sparse than when all predictors are of interest, which is the case in the simulated toy examples. Indeed, on such datasets randomization will lead to neglect relevant information, while in sparse cases it will improve the chances to select the most relevant variables.

On these toy examples, best performances in terms of total and local AUC and of standard type error are achieved by the $RF1_{SVM}$ heuristic and interestingly, the algorithms reaching the best performance in terms of low variance are the ones with high randomization ($RF3_{CART}$, $RF4_{CART}$, $RF5_{CART}$, $RF6_{CART}$). Yet, given the weaker ranking performance of these highly randomized versions (in comparison to the others), we discard the following heuristics $RF2_{CART}$, $RF2_{SVM}$, $RF4_{CART}$ and $RF6_{CART}$ from the subsequent comparisons.

In the following series of experiments, we focus on the eight remaining heuristics that we compare on real benchmark data sets, gathered from the UCI repository. This time, model selection is conducted via the pruning procedure introduced in 3.1. Therefore, the regularization coefficient $\lambda$ has been calibrated through a nested stratified 5-fold cross validation procedure. As previously, a global overview of the performances overall the benchmark datasets is provided by Figure 5. Visual displays are also detailed for some of the datasets on Figure 6 below. Again, the histograms represent performance and standard type error ($AUC, \hat{\sigma}^2$), as well as Local AUC measure $AUC@u\%$ for a proportion $u = 20\%$ of best instances.

The quantitative results can be found in Tables 5 and 6, gathering full and local AUC for



**Fig. 5** Average ranks of TREERANK variants on the 10 benchmark datasets for several performance indicators.

a. Dataset Congressional Vote

b. Dataset Australian Credit

c. Dataset German Credit

d. Dataset Autos MPG

e. Dataset Ionosphere

f. Dataset Heart Disease

**Fig. 6** Comparison of TREERANK variants on benchmark datasets.

proportions $u = 20, 10, 5\%$ of best instances indicated together with their associated standard type errors, given in parentheses, as well as hit ratios for proportions 10% and 20%, and the average precision. All these values are computed through a stratified 10-fold cross validation, as described in 4.3.

The global picture displayed on Figure 5 clearly shows two dynamics. When considering the indicators related to learning performances, we notice that randomized procedure are among the most efficient. Moreover, different behaviours appear when considering the ro-

bustness indicators. Indeed, the graph shows that the algorithms with highest AUC, average precision and hit ratios are also the ones having highest standard deviations, and conversely. This is due to the fact that these procedures are randomized. Indeed, the only procedure presenting both good learning performances and robustness is $\text{Bagg}_{\text{SVM}}$.

Besides, as noticed on the simulated toy examples, we can see on the graphs of Figure 6 that in most cases, making use of SVM classifiers as splitting rule leads to more efficient ranking rules, in terms of both full and Lo-

cal AUC, but also in terms of stability (lower standard deviation). As explained previously, SVM based classifier allow to capture the complex geometry of the regression level sets in a more flexible manner than CART based decision rules. Yet, when focusing on the *Heart Disease* dataset for instance, we can observe that Bagg$_{\text{SVM}}$ and RF1$_{\text{SVM}}$ perform worse on the 20% best instances than their CART-based counterparts Bagg$_{\text{CART}}$ and RF1$_{\text{CART}}$. This observation means that in this particular case, the topology of the first level sets (corresponding to the 20% of best instances) is better approximated by linear partitions than by nonlinear Gaussian partitions. Indeed, the geometry of the level sets can be different from one dataset to another, but also from one level set to another. This example emphasizes the advantage of using different splitting rules in the TREE-RANK procedure, its performance being intrinsically related to the capacity of the latter to recover the geometry of the level sets. Moreover, when comparing TRK$_{\text{CART}}$ and TRK$_{\text{SVM}}$ with their *resampled* versions, the positive impact of resampling and aggregation procedures on the performances and the variance of the output rules clearly appears.

Eventually, the main difference with the previous experiments lies in the impact of feature randomization. Even if Bagg$_{\text{SVM}}$ achieves very good performances on several datasets, the best results are generally provided by resampled and randomized versions. Three heuristics perform particularly better than the others: Bagg$_{\text{SVM}}$ and RF1$_{\text{SVM}}$, which often get similar results, and RF3$_{\text{CART}}$. For the following series of experiments we will retain only two (RF1$_{\text{SVM}}$ and RF3$_{\text{CART}}$), which we will compare to several competitors.

TREERANK *and some competitors.* In this series of experiments, both RF1$_{\text{SVM}}$ and RF3$_{\text{CART}}$ heuristics are compared to six scoring algorithms: RANKBOOST, RANKLS, RANKSVM, KLR, AD-ABOOST and P-NORM PUSH. We first consider the performances of these heuristics on the three toy examples previously described, before considering the benchmark datasets from the UCI repository. As before, for the experiments based on simulated datasets the maximum depths of the (sub-) trees has been settled *a priori* respectively to $d_1 = 8$ and $d_2 = 6$, while on the UCI benchmark datasets model selection is based on



**Fig. 7** Average ranks of TREERANK variants and some competitors on the toy examples for several performance indicators.

a nested stratified 5-fold cross validation procedure. The quantitative results obtained on synthetic data are summarized in Tables 7, 8 and 9. They are also displayed on Figure 8 and a global overview of the performances of each algorithm is displayed on Figure 7.

As in the previous series of experiments, we notice that generally the algorithms with highest performances are also the less robust (see Figure 7). Moreover, we can see on Figure 8 that the performances of all heuristics in terms of global and local AUC are very similar on these data. Yet, we can observe that RANKLS seems to perform a little better on the *GaussHard20d* data set. In the other cases the leadership goes to a version of RANKING FOREST with an SVM classifier used as splitting rule. It is also interesting to notice that the instability results are much more variable than in previous experiments. Indeed, there is a big difference between the algorithms and from a data set to another. ADABOOST and RF3$_{\text{CART}}$ heuristics appear to be more stable than the others. We can also underline that for several algorithms, variance increases with the complexity of the ranking problem (from *GaussEasy20d* data set to *GaussHard20d*). It is the case in particular for RANKBOOST, RANKLS, RANKSVM and KLR.

Yet, differences between the several algorithms appear more clearly on the real datasets. Performance and robustness results are summarized in Tables 10 and 11, detailed graphs per dataset are displayed on Figure 10 and a global comparison of the algorithms is given by Figure 9.

a. Dataset GaussEasy20D



b. Dataset GaussMed20D



c. Dataset GaussHard20D

**Fig. 8** Comparison of TREERANK and some competitors on simulated datasets.

We clearly see on Figure 9 that both RANK-ING FOREST algorithms, either based on the LEAFRANK procedure or on a SVM classifier, perform mostly better than other methods. This is particularly true regarding learning performances, while both methods present lower robustness. When going into more detail, we notice three exceptions: Australian, German and Japanese credit datsets, on which RANKBOOST

achieves better performances than TREERANK (see Figure 10). However, TREERANK heuristics clearly provide better results on the best instances for all the considered datasets. This is related to the fact that TREERANK recursively optimizes the AUC and thus estimates the optimal ROC curve, while all other methods optimize the AUC criterion in a global manner.

## 5 Discussion

The numerical experiments performed in the present paper show that the three methodologies (Ranking trees, Pairwise SVM, Pairwise Boosting) are very competitive for AUC maximization. It is shown that the ranking trees heuristics take the lead when it comes to focusing on best instances. Theoretical evidence supports this experimental fact. Indeed, in their very construction, ranking trees derived from the TREERANK/LEAFRANK algorithm actually optimize the AUC locally. The present paper illustrates the theoretical result introduced in



**Fig. 9** Average ranks of TREERANK variants and some competitors on the 10 benchmark datasets for several performance indicators.

a. Dataset Congressional Vote



b. Dataset Australian Credit



c. Dataset German Credit



d. Dataset Autos MPG



e. Dataset Ionosphere



f. Dataset Heart Disease

**Fig. 10** Comparison of TREERANK and some competitors on benchmark datasets.

([1]) which states that, by using ranking trees, convergence to the optimal ROC curve is performed in a stronger sense than the AUC. Indeed, AUC optimality corresponds to convergence in $L_1$-norm of the corresponding ROC curve, while it can be shown that the recursive construction of ranking trees guarantees convergence in $L_\infty$-norm. This property explains the improvement obtained with our method with respect to local ranking performance measures.

| AUC AUC@20% AUC@10% AUC@5% | *GaussEasy20d* | *GaussMed20d* | *GaussHard20d* |
|---|---|---|---|
| TRK$_{\text{CART}}$ | 0.626 ($\pm$0.015) 0.243 ($\pm$0.011) 0.130 ($\pm$0.010) 0.067 ($\pm$0.006) | 0.625 ($\pm$0.017) 0.244 ($\pm$0.013) 0.132 ($\pm$0.012) 0.068 ($\pm$0.006) | 0.539 ($\pm$0.018) 0.195 ($\pm$0.015) 0.102 ($\pm$0.009) 0.053 ($\pm$0.006) |
| TRK$_{\text{SVM}}$ | 0.752 ($\pm$0.009) 0.292 ($\pm$0.011) 0.150 ($\pm$0.008) 0.075 ($\pm$0.004) | 0.721 ($\pm$0.011) 0.282 ($\pm$0.015) 0.145 ($\pm$0.009) 0.073 ($\pm$0.005) | 0.562 ($\pm$0.017) 0.201 ($\pm$0.015) 0.103 ($\pm$0.010) 0.052 ($\pm$0.006) |
| Bagg$_{\text{CART}}$ | 0.722 ($\pm$0.010) 0.300 ($\pm$0.007) 0.163 ($\pm$0.006) 0.085 ($\pm$0.005) | 0.714 ($\pm$0.013) 0.297 ($\pm$0.009) 0.163 ($\pm$0.006) 0.086 ($\pm$0.004) | 0.566 ($\pm$0.014) 0.214 ($\pm$0.011) 0.113 ($\pm$0.009) 0.060 ($\pm$0.007) |
| Bagg$_{\text{SVM}}$ | 0.772 ($\pm$0.008) 0.298 ($\pm$0.010) 0.156 ($\pm$0.007) 0.080 ($\pm$0.005) | 0.739 ($\pm$0.009) 0.287 ($\pm$0.012) 0.151 ($\pm$0.008) 0.078 ($\pm$0.005) | 0.576 ($\pm$0.014) 0.207 ($\pm$0.013) 0.107 ($\pm$0.010) 0.055 ($\pm$0.007) |
| RF1$_{\text{CART}}$ | 0.747 ($\pm$0.010) 0.308 ($\pm$0.009) 0.167 ($\pm$0.005) 0.087 ($\pm$0.004) | 0.737 ($\pm$0.012) 0.308 ($\pm$0.009) 0.165 ($\pm$0.006) 0.087 ($\pm$0.004) | 0.569 ($\pm$0.013) 0.211 ($\pm$0.011) 0.113 ($\pm$0.008) ;0.060 ($\pm$0.006) |
| RF1$_{\text{SVM}}$ | **0.776** ($\pm$**0.006**) **0.319** ($\pm$**0.006**) **0.172** ($\pm$**0.004**) **0.089** ($\pm$**0.003**) | **0.760** ($\pm$**0.006**) **0.316** ($\pm$**0.005**) **0.171** ($\pm$**0.003**) **0.089** ($\pm$**0.003**) | **0.582** ($\pm$**0.013**) **0.216** ($\pm$**0.007**) **0.115** ($\pm$**0.007**) **0.061** ($\pm$**0.006**) |
| RF2$_{\text{CART}}$ | 0.718 ($\pm$0.011) 0.294 ($\pm$0.011) 0.162 ($\pm$0.007) 0.086 ($\pm$0.003) | 0.701 ($\pm$0.014) 0.289 ($\pm$0.011) 0.157 ($\pm$0.008) 0.083 ($\pm$0.004) | 0.557 ($\pm$0.012) 0.205 ($\pm$0.009) 0.109 ($\pm$0.007) 0.056 ($\pm$0.007) |
| RF2$_{\text{SVM}}$ | 0.697 ($\pm$0.015) 0.280 ($\pm$0.012) 0.152 ($\pm$0.007) 0.081 ($\pm$0.005) | 0.675 ($\pm$0.016) 0.271 ($\pm$0.011) 0.145 ($\pm$0.008) 0.077 ($\pm$0.005) | 0.539 ($\pm$0.016) 0.197 ($\pm$0.009) 0.105 ($\pm$0.008) 0.055 ($\pm$0.007) |
| RF3$_{\text{CART}}$ | 0.718 ($\pm$0.013) 0.289 ($\pm$0.010) 0.156 ($\pm$0.008) 0.081 ($\pm$0.005) | 0.698 ($\pm$0.016) 0.283 ($\pm$0.013) 0.155 ($\pm$0.006) 0.081 ($\pm$0.005) | 0.551 ($\pm$0.019) 0.203 ($\pm$0.012) 0.108 ($\pm$0.009) 0.057 ($\pm$0.007) |
| RF4$_{\text{CART}}$ | 0.693 ($\pm$0.017) 0.278 ($\pm$0.013) 0.151 ($\pm$0.009) 0.080 ($\pm$0.006) | 0.670 ($\pm$0.015) 0.267 ($\pm$0.013) 0.144 ($\pm$0.009) 0.078 ($\pm$0.005) | 0.540 ($\pm$0.021) 0.200 ($\pm$0.013) 0.106 ($\pm$0.008) 0.055 ($\pm$0.005) |
| RF5$_{\text{CART}}$ | 0.684 ($\pm$0.021) 0.269 ($\pm$0.014) 0.148 ($\pm$0.008) 0.078 ($\pm$0.005) | 0.662 ($\pm$0.028) 0.265 ($\pm$0.017) 0.144 ($\pm$0.010) 0.076 ($\pm$0.006) | 0.535 ($\pm$0.028) 0.195 ($\pm$0.016) 0.106 ($\pm$0.011) 0.054 ($\pm$0.008) |
| RF6$_{\text{CART}}$ | 0.622 ($\pm$0.037) 0.244 ($\pm$0.019) 0.131 ($\pm$0.012) 0.068 ($\pm$0.007) | 0.592 ($\pm$0.029) 0.225 ($\pm$0.018) 0.123 ($\pm$0.014) 0.064 ($\pm$0.010) | 0.519 ($\pm$0.028) 0.189 ($\pm$0.017) 0.099 ($\pm$0.013) 0.053 ($\pm$0.008) |

**Table 2** Comparison of 12 variants of the TREERANK algorithm on simulated data. Performance is measured through the AUC and Local AUC for three different values $u = 20, 10, 5\%$, given in that order. Standard type error $\hat{\sigma}^2$ for these values is indicated in parentheses.

| HR@10% ; HR@20% AP | GaussEasy20d | GaussMed20d | GaussHard20d |
|---|---|---|---|
| TRK$_{CART}$ | 68%; 66% 0.59 | 68%; 65% 0.58 | 54%; 55% 0.53 |
| TRK$_{SVM}$ | 79%; 78% 0.67 | 75%; 74% 0.64 | 54%; 56% 0.54 |
| Bagg$_{CART}$ | 84%; 79% 0.68 | 82%; 77% 0.66 | 60%; 59% 0.56 |
| Bagg$_{SVM}$ | 81%; 79% 0.69 | 77%; 75% 0.66 | 57%; 57% 0.55 |
| RF1$_{CART}$ | 86%; 80% 0.69 | 84%; 79% 0.68 | 60%; 58% 0.56 |
| RF1$_{SVM}$ | **88**%; **83**% **0.71** | **86**%; **81**% **0.69** | **61**%; **59**% **0.56** |
| RF2$_{CART}$ | 84%; 77% 0.67 | 79%; 75% 0.65 | 58%; 57% 0.55 |
| RF2$_{SVM}$ | 79%; 74% 0.66 | 74%; 71% 0.63 | 56%; 55% 0.53 |
| RF3$_{CART}$ | 81%; 76% 0.67 | 79%; 73% 0.65 | 57%; 56% 0.54 |
| RF4$_{CART}$ | 78%; 74% 0.65 | 74%; 70% 0.62 | 57%; 55% 0.53 |
| RF5$_{CART}$ | 77%; 72% 0.64 | 74%; 69% 0.62 | 56%; 54% 0.53 |
| RF6$_{CART}$ | 69%; 66% 0.60 | 64%; 60% 0.57 | 53%; 53% 0.52 |

**Table 3** Comparison of 12 variants of the TREERANK algorithm on simulated data. Performance is measured through the hit ratios for proportions 10% and 20%, respectively given in that order, and the average precision indicated below.

| $\widehat{\text{Instab}}$ | GaussEasy20d | GaussMed20d | GaussHard20d |
|---|---|---|---|
| TRK$_{CART}$ | 137.5 | 141.7 | 149.5 |
| TRK$_{SVM}$ | 56.5 | 64.5 | 108.2 |
| Bagg$_{CART}$ | 15.3 | 16.3 | 16.3 |
| Bagg$_{SVM}$ | 24.4 | 27.8 | 41.3 |
| RF1$_{CART}$ | 6.5 | 7.2 | 7.9 |
| RF1$_{SVM}$ | 5.5 | 6.4 | 9.6 |
| RF2$_{CART}$ | 4.5 | 4.8 | 5.7 |
| RF2$_{SVM}$ | 5.6 | 6.1 | 7.2 |
| RF3$_{CART}$ | 3.5 | 3.5 | 3.6 |
| RF4$_{CART}$ | 3.1 | 3.3 | 3.2 |
| RF5$_{CART}$ | 3.4 | 3.3 | 3.3 |
| RF6$_{CART}$ | **3.0** | **3.2** | **3.1** |

**Table 4** Comparison of 12 variants of the TREERANK algorithm on simulated data in terms of instability, computed with $B = 20$ bootstrap samples.

| AUC<br>AUC@20%<br>AUC@10%<br>AUC@5% | $\text{TRK}_{\text{CART}}$ | $\text{TRK}_{\text{SVM}}$ | $\text{Bagg}_{\text{CART}}$ | $\text{Bagg}_{\text{SVM}}$ | $\text{RF1}_{\text{CART}}$ | $\text{RF1}_{\text{SVM}}$ | $\text{RF3}_{\text{CART}}$ | $\text{RF5}_{\text{CART}}$ |
|---|---|---|---|---|---|---|---|---|
| *Congressional Vote* | 0.749 ($\pm$0.096)<br>0.261 ($\pm$0.060)<br>0.128 ($\pm$0.052)<br>0.069 ($\pm$0.046) | 0.991 ($\pm$0.024)<br>0.595 ($\pm$0.105)<br>0.186 ($\pm$0.004)<br>0.093 ($\pm$0.002) | 0.905 ($\pm$0.06)<br>0.359 ($\pm$0.049)<br>0.232 ($\pm$0.075)<br>0.11 ($\pm$0.031) | 0.991 ($\pm$0.026)<br>0.650 ($\pm$0.112)<br>0.186 ($\pm$0.004)<br>0.093 ($\pm$0.002) | 0.985 ($\pm$0.028)<br>0.473 ($\pm$0.061)<br>0.217 ($\pm$0.064)<br>0.094 ($\pm$0) | 0.987 ($\pm$0.024)<br>0.477 ($\pm$0.054)<br>0.201 ($\pm$0.042)<br>0.094 ($\pm$0) | **0.992 ($\pm$0.015)**<br>**0.495 ($\pm$0.059)**<br>0.201 ($\pm$0.042)<br>0.094 ($\pm$0) | 0.988 ($\pm$0.021)<br>0.474 ($\pm$0.053)<br>0.204 ($\pm$0.054)<br>0.094 ($\pm$0) |
| *Australian Credit* | 0.577 ($\pm$0.057)<br>0.257 ($\pm$0.050)<br>0.134 ($\pm$0.023)<br>0.068 ($\pm$0.011) | 0.921 ($\pm$0.032)<br>0.414 ($\pm$0.024)<br>0.256 ($\pm$0.082)<br>0.115 ($\pm$0.028) | 0.852 ($\pm$0.038)<br>0.419 ($\pm$0.014)<br>0.235 ($\pm$0.031)<br>0.141 ($\pm$0.041) | 0.925 ($\pm$0.034)<br>0.423 ($\pm$0.016)<br>0.252 ($\pm$0.040)<br>0.115 ($\pm$0.012) | 0.882 ($\pm$0.026)<br>0.771 ($\pm$0.059)<br>0.663 ($\pm$0.075)<br>0.543 ($\pm$0.081) | 0.926 ($\pm$0.031)<br>0.425 ($\pm$0.012)<br>0.248 ($\pm$0.039)<br>0.111 ($\pm$0.002) | 0.934 ($\pm$0.031)<br>0.429 ($\pm$0.014)<br>0.273 ($\pm$0.059)<br>0.122 ($\pm$0.031) | **0.939 ($\pm$0.031)**<br>**0.429 ($\pm$0.014)**<br>0.278 ($\pm$0.061)<br>0.129 ($\pm$0.038) |
| *German Credit* | 0.711 ($\pm$0.071)<br>0.237 ($\pm$0.031)<br>0.123 ($\pm$0.019)<br>0.061 ($\pm$0.014) | 0.771 ($\pm$0.034)<br>0.254 ($\pm$0.017)<br>0.133 ($\pm$0.018)<br>0.065 ($\pm$0.004) | 0.737 ($\pm$0.045)<br>0.245 ($\pm$0.022)<br>0.128 ($\pm$0.020)<br>0.063 ($\pm$0.008) | **0.794 ($\pm$0.024)**<br>**0.269 ($\pm$0.009)**<br>0.146 ($\pm$0.016)<br>0.098 ($\pm$0.033) | 0.771 ($\pm$0.031)<br>0.264 ($\pm$0.016)<br>0.146 ($\pm$0.021)<br>0.075 ($\pm$0.014) | 0.793 ($\pm$0.022)<br>0.266 ($\pm$0.012)<br>0.148 ($\pm$0.020)<br>0.081 ($\pm$0.021) | 0.769 ($\pm$0.029)<br>0.260 ($\pm$0.020)<br>0.144 ($\pm$0.022)<br>0.083 ($\pm$0.029) | 0.775 ($\pm$0.044)<br>0.263 ($\pm$0.017)<br>0.147 ($\pm$0.027)<br>0.077 ($\pm$0.015) |
| *Japanese Credit* | 0.789 ($\pm$0.053)<br>0.373 ($\pm$0.043)<br>0.220 ($\pm$0.059)<br>0.106 ($\pm$0.024) | 0.924 ($\pm$0.049)<br>0.423 ($\pm$0.050)<br>0.232 ($\pm$0.053)<br>0.105 ($\pm$0.005) | 0.839 ($\pm$0.048)<br>0.406 ($\pm$0.024)<br>0.250 ($\pm$0.047)<br>0.131 ($\pm$0.035) | 0.926 ($\pm$0.044)<br>0.416 ($\pm$0.026)<br>0.245 ($\pm$0.052)<br>0.127 ($\pm$0.030) | 0.887 ($\pm$0.038)<br>0.416 ($\pm$0.015)<br>0.256 ($\pm$0.063)<br>0.125 ($\pm$0.033) | 0.923 ($\pm$0.038)<br>0.414 ($\pm$0.020)<br>0.256 ($\pm$0.068)<br>0.125 ($\pm$0.032) | **0.931 ($\pm$0.042)**<br>**0.417 ($\pm$0.023)**<br>0.260 ($\pm$0.073)<br>0.141 ($\pm$0.047) | **0.930 ($\pm$0.043)**<br>**0.418 ($\pm$0.027)**<br>0.274 ($\pm$0.081)<br>0.123 ($\pm$0.035) |
| *Autos MPG* | 0.927 ($\pm$0.049)<br>0.420 ($\pm$0.011)<br>0.229 ($\pm$0.081)<br>0.092 ($\pm$0.007) | 0.961 ($\pm$0.041)<br>0.462 ($\pm$0.144)<br>0.184 ($\pm$0.010)<br>0.092 ($\pm$0.005) | 0.967 ($\pm$0.022)<br>0.523 ($\pm$0.112)<br>0.190 ($\pm$0.002)<br>0.094 ($\pm$0.004) | **0.982 ($\pm$0.017)**<br>**0.584 ($\pm$0.109)**<br>0.190 ($\pm$0)<br>0.095 ($\pm$0) | 0.978 ($\pm$0.014)<br>0.418 ($\pm$0.039)<br>0.244 ($\pm$0.087)<br>0.095 ($\pm$0) | 0.978 ($\pm$0.016)<br>0.433 ($\pm$0.058)<br>0.238 ($\pm$0.079)<br>0.095 ($\pm$0) | 0.972 ($\pm$0.016)<br>0.421 ($\pm$0.016)<br>0.240 ($\pm$0.08)<br>0.095 ($\pm$0) | 0.964 ($\pm$0.021)<br>0.418 ($\pm$0.046)<br>0.236 ($\pm$0.076)<br>0.095 ($\pm$0) |
| *Ionosphere* | 0.926 ($\pm$0.042)<br>0.360 ($\pm$0.127)<br>0.146 ($\pm$0.007)<br>0.074 ($\pm$0.003) | 0.943 ($\pm$0.049)<br>0.352 ($\pm$0.129)<br>0.149 ($\pm$0.008)<br>0.074 ($\pm$0.004) | 0.966 ($\pm$0.026)<br>0.408 ($\pm$0.116)<br>0.181 ($\pm$0.052)<br>0.078 ($\pm$0) | 0.985 ($\pm$0.014)<br>0.477 ($\pm$0.072)<br>0.156 ($\pm$0)<br>0.078 ($\pm$0) | 0.971 ($\pm$0.026)<br>0.411 ($\pm$0.099)<br>0.179 ($\pm$0.048)<br>0.078 ($\pm$0) | **0.990 ($\pm$0.013)**<br>**0.494 ($\pm$0.062)**<br>0.156 ($\pm$0)<br>0.078 ($\pm$0) | 0.968 ($\pm$0.032)<br>0.392 ($\pm$0.086)<br>0.178 ($\pm$0.044)<br>0.078 ($\pm$0) | 0.966 ($\pm$0.032)<br>0.391 ($\pm$0.032)<br>0.179 ($\pm$0.038)<br>0.081 ($\pm$0.008) |
| *Breast Cancer Diagnosis* | 0.958 ($\pm$0.02)<br>0.503 ($\pm$0.015)<br>0.314 ($\pm$0.059)<br>0.131 ($\pm$0.004) | 0.989 ($\pm$0.008)<br>0.565 ($\pm$0.063)<br>0.341 ($\pm$0.102)<br>0.133 ($\pm$0.003) | 0.986 ($\pm$0.016)<br>0.578 ($\pm$0.036)<br>0.319 ($\pm$0.106)<br>0.134 ($\pm$0) | 0.990 ($\pm$0.015)<br>**0.680 ($\pm$0.071)**<br>0.290 ($\pm$0.070)<br>0.134 ($\pm$0) | 0.990 ($\pm$0.009)<br>0.548 ($\pm$0.025)<br>0.403 ($\pm$0.121)<br>0.134 ($\pm$0) | **0.994 ($\pm$0.008)**<br>0.595 ($\pm$0.047)<br>0.317 ($\pm$0.103)<br>0.134 ($\pm$0) | 0.988 ($\pm$0.007)<br>0.528 ($\pm$0.005)<br>0.453 ($\pm$0.036)<br>0.135 ($\pm$0) | 0.990 ($\pm$0.007)<br>0.527 ($\pm$0.006)<br>0.444 ($\pm$0.037)<br>0.134 ($\pm$0) |
| *Breast Cancer Original* | 0.984 ($\pm$0.017)<br>0.568 ($\pm$0.024)<br>0.375 ($\pm$0.101)<br>0.154 ($\pm$0.039) | 0.995 ($\pm$0.005)<br>0.566 ($\pm$0.027)<br>0.364 ($\pm$0.107)<br>0.139 ($\pm$0.006) | 0.990 ($\pm$0.009)<br>0.565 ($\pm$0.014)<br>0.441 ($\pm$0.106)<br>0.157 ($\pm$0.044) | 0.994 ($\pm$0.007)<br>**0.570 ($\pm$0.027)**<br>0.399 ($\pm$0.133)<br>0.157 ($\pm$0.043) | **0.995 ($\pm$0.005)**<br>0.558 ($\pm$0.009)<br>0.430 ($\pm$0.085)<br>0.148 ($\pm$0.016) | 0.995 ($\pm$0.006)<br>0.559 ($\pm$0.010)<br>0.442 ($\pm$0.076)<br>0.146 ($\pm$0.010) | 0.993 ($\pm$0.006)<br>0.556 ($\pm$0.012)<br>0.407 ($\pm$0.078)<br>0.142 ($\pm$0.002) | 0.993 ($\pm$0.006)<br>0.555 ($\pm$0.010)<br>0.409 ($\pm$0.087)<br>0.156 ($\pm$0.042) |
| *Heart Disease* | 0.794 ($\pm$0.061)<br>0.376 ($\pm$0.054)<br>0.237 ($\pm$0.092)<br>0.109 ($\pm$0.037) | 0.868 ($\pm$0.048)<br>0.401 ($\pm$0.057)<br>0.238 ($\pm$0.090)<br>0.101 ($\pm$0.011) | 0.798 ($\pm$0.074)<br>0.372 ($\pm$0.047)<br>0.221 ($\pm$0.056)<br>0.139 ($\pm$0.058) | 0.908 ($\pm$0.049)<br>0.407 ($\pm$0.050)<br>0.278 ($\pm$0.099)<br>0.114 ($\pm$0.030) | 0.869 ($\pm$0.067)<br>0.397 ($\pm$0.042)<br>0.261 ($\pm$0.064)<br>0.118 ($\pm$0.033) | **0.917 ($\pm$0.032)**<br>**0.416 ($\pm$0.027)**<br>0.273 ($\pm$0.070)<br>0.118 ($\pm$0.017) | 0.911 ($\pm$0.034)<br>0.415 ($\pm$0.026)<br>0.272 ($\pm$0.070)<br>0.121 ($\pm$0.031) | 0.907 ($\pm$0.048)<br>0.408 ($\pm$0.033)<br>0.289 ($\pm$0.076)<br>0.120 ($\pm$0.032) |
| *Hepatitis* | 0.771 ($\pm$0.171)<br>0.405 ($\pm$0.146)<br>0.25 ($\pm$0.139)<br>0.174 ($\pm$0.142) | 0.813 ($\pm$0.154)<br>0.547 ($\pm$0.202)<br>0.306 ($\pm$0.129)<br>0.188 ($\pm$0.137) | 0.813 ($\pm$0.132)<br>0.470 ($\pm$0.132)<br>0.279 ($\pm$0.172)<br>0.202 ($\pm$0.177) | **0.872 ($\pm$0.14)**<br>**0.629 ($\pm$0.213)**<br>0.412 ($\pm$0.131)<br>0.339 ($\pm$0.157) | 0.850 ($\pm$0.141)<br>0.551 ($\pm$0.24)<br>0.327 ($\pm$0.178)<br>0.251 ($\pm$0.196) | 0.864 ($\pm$0.139)<br>0.572 ($\pm$0.240)<br>0.413 ($\pm$0.138)<br>0.269 ($\pm$0.190) | 0.868 ($\pm$0.162)<br>0.625 ($\pm$0.230)<br>0.346 ($\pm$0.162)<br>0.342 ($\pm$0.154) | 0.842 ($\pm$0.17)<br>0.562 ($\pm$0.233)<br>0.277 ($\pm$0.187)<br>0.277 ($\pm$0.187) |

**Table 5** Comparison of eight variants of the TREERANK algorithm on benchmark data. Performance is measured through AUC and Local AUC for three different values $u = 20, 10, 5\%$, given in that order. Standard type error $\hat{\sigma}^2$ for these values is indicated in parentheses.

| HR@10% ; HR@20% AP | TRK$_{CART}$ | TRK$_{SVM}$ | Bagg$_{CART}$ | Bagg$_{SVM}$ | RF1$_{CART}$ | RF1$_{SVM}$ | RF3$_{CART}$ | RF5$_{CART}$ |
|---|---|---|---|---|---|---|---|---|
| *Congressional Vote* | 92%; 92% 0.75 | 100%; 98% 0.80 | 100%; 96% 0.80 | 100%; 98% 0.82 | 100%; 100% 0.84 | 100%; 100% 0.84 | **100%; 100%** **0.85** | 100%; 100% 0.84 |
| *Australian Credit* | 68%; 59% 0.48 | 100%; 100% 0.72 | 100%; 100% 0.71 | 100%; 100% 0.75 | 100%; 100% 0.73 | 100%; 100% 0.75 | 100%; 100% 0.76 | **100%; 100%** **0.76** |
| *German Credit* | 87%; 84% 0.80 | 94%; 91% 0.84 | 90%; 88% 0.82 | **99%; 96%** **0.87** | 96%; 93% 0.85 | 98%; 94% 0.86 | 99%; 92% 0.86 | 99%; 93% 0.86 |
| *Japanese Credit* | 78%; 83% 0.64 | 97%; 95% 0.68 | 98%; 92% 0.70 | 99%; 95% 0.71 | 98%; 97% 0.70 | 95%; 95% 0.71 | **97%; 95%** **0.71** | **97%; 95%** **0.72** |
| *Autos MPG* | 89%; 92% 0.74 | 87%; 92% 0.72 | 98%; 99% 0.78 | 100%; 100% 0.77 | 100%; 100% 0.79 | **100%; 100%** **0.80** | 100%; 100% 0.79 | **100%; 100%** **0.80** |
| *Ionosphere* | 100%; 100% 0.85 | 100%; 97% 0.87 | 100%; 99% 0.91 | 100%; 100% 0.93 | 100%; 100% 0.92 | **100%; 100%** **0.94** | 100%; 99% 0.92 | 100%; 99% 0.92 |
| *Breast Cancer Diagnosis* | 98%; 97% 0.57 | 98%; 99% 0.60 | 100%; 100% 0.69 | 100%; 100% 0.67 | 100%; 100% 0.69 | **100%; 100%** **0.70** | 100%; 100% 0.70 | 100%; 100% 0.70 |
| *Breast Cancer Original* | 100%; 99% 0.52 | 99%; 99% 0.50 | 100%; 99% 0.60 | 99%; 99% 0.55 | **100%; 99%** **0.64** | 100%; 99% 0.65 | 99%; 99% 0.64 | 100%; 100% 0.64 |
| *Heart Disease* | 71%; 72% 0.59 | 90%; 89% 0.67 | 93%; 86% 0.66 | 90%; 90% 0.74 | 93%; 94% 0.69 | **97%; 94%** **0.73** | 93%; 96% 0.75 | 93%; 94% 0.72 |
| *Hepatitis* | 75%; 40% 0.25 | 65%; 47% 0.31 | 45%; 35% 0.29 | **75%; 47%** **0.41** | 45%; 43% 0.31 | 65%; 40% 0.33 | **85%; 42%** **0.43** | 65%; 42% 0.31 |

**Table 6** Comparison of eight variants of the TREERANK algorithm on benchmark data. Performance is measured through the hit ratios for proportions 10% and 20%, given in that order respectively, and the average precision given below.

| AUC AUC@20% AUC@10% AUC@5% | *GaussEasy20d* | *GaussMed20d* | *GaussHard20d* |
|---|---|---|---|
| **RF1**$_{\text{SVM}}$ | **0.776** ($\pm$**0.006**) | **0.760** ($\pm$**0.006**) | **0.582** ($\pm$**0.013**) |
| | **0.319** ($\pm$**0.006**) | **0.316** ($\pm$**0.005**) | **0.216** ($\pm$**0.007**) |
| | 0.172 ($\pm$0.004) | 0.171 ($\pm$0.003) | 0.115 ($\pm$0.007) |
| | 0.089 ($\pm$0.003) | 0.089 ($\pm$0.003) | 0.061 ($\pm$0.006) |
| **RF3**$_{\text{CART}}$ | 0.718 ($\pm$0.013) | 0.698 ($\pm$0.016) | 0.551 ($\pm$0.019) |
| | 0.289 ($\pm$0.010) | 0.283 ($\pm$0.013) | 0.203 ($\pm$0.012) |
| | 0.156 ($\pm$0.008) | 0.155 ($\pm$0.006) | 0.108 ($\pm$0.009) |
| | 0.081 ($\pm$0.005) | 0.081 ($\pm$0.005) | 0.057 ($\pm$0.007) |
| AdaBoost | 0.716 ($\pm$0.009) | 0.714 ($\pm$0.010) | 0.565 ($\pm$0.015) |
| | 0.286 ($\pm$0.007) | 0.289 ($\pm$0.008) | 0.209 ($\pm$0.010) |
| | 0.152 ($\pm$0.005) | 0.158 ($\pm$0.006) | 0.111 ($\pm$0.008) |
| | 0.080 ($\pm$0.005) | 0.083 ($\pm$0.004) | 0.057 ($\pm$0.005) |
| RankBoost | **0.754** ($\pm$**0.007**) | **0.746** ($\pm$**0.008**) | **0.577** ($\pm$**0.015**) |
| | **0.306** ($\pm$**0.009**) | **0.307** ($\pm$**0.009**) | **0.577** ($\pm$**0.015**) |
| | 0.163 ($\pm$0.006) | 0.166 ($\pm$0.004) | 0.111 ($\pm$0.007) |
| | 0.087 ($\pm$0.006) | 0.086 ($\pm$0.004) | 0.058 ($\pm$0.005) |
| RankSVM | 0.740 ($\pm$0.004) | 0.737 ($\pm$0.007) | **0.578** ($\pm$**0.011**) |
| | 0.282 ($\pm$0.005) | 0.288 ($\pm$0.007) | **0.213** ($\pm$**0.010**) |
| | 0.147 ($\pm$0.003) | 0.153 ($\pm$0.005) | 0.112 ($\pm$0.008) |
| | 0.075 ($\pm$0.003) | 0.081 ($\pm$0.004) | 0.056 ($\pm$0.005) |
| RankLS | 0.742 ($\pm$0.004) | 0.719 ($\pm$0.015) | **0.581** ($\pm$**0.010**) |
| | 0.283 ($\pm$0.005) | 0.278 ($\pm$0.008) | **0.216** ($\pm$**0.008**) |
| | 0.147 ($\pm$0.003) | 0.151 ($\pm$0.006) | 0.112 ($\pm$0.007) |
| | 0.075 ($\pm$0.003) | 0.080 ($\pm$0.005) | 0.057 ($\pm$0.005) |
| P − norm push | 0.700 ($\pm$0.017) | 0.701 ($\pm$0.014) | 0.566 ($\pm$0.016) |
| | 0.271 ($\pm$0.010) | 0.272 ($\pm$0.009) | 0.208 ($\pm$0.011) |
| | 0.142 ($\pm$0.006) | 0.146 ($\pm$0.006) | 0.109 ($\pm$0.008) |
| | 0.071 ($\pm$0.004) | 0.076 ($\pm$0.005) | 0.056 ($\pm$0.005) |
| KLR | 0.742 ($\pm$0.002) | 0.740 ($\pm$0.005) | **0.582** ($\pm$**0.011**) |
| | 0.284 ($\pm$0.004) | 0.286 ($\pm$0.006) | **0.215** ($\pm$**0.010**) |
| | 0.147 ($\pm$0.004) | 0.155 ($\pm$0.004) | 0.112 ($\pm$0.008) |
| | 0.075 ($\pm$0.003) | 0.082 ($\pm$0.005) | 0.057 ($\pm$0.005) |

**Table 7** Comparison of TREERANK and some competitors on simulated data. Total AUC and Local AUC for three different values $u = 20, 10, 5\%$ are given in that order. Standard type error $\hat{\sigma}^2$ for these values is indicated in parentheses.

| HR@10% ; HR@20% AP | *GaussEasy20d* | *GaussMed20d* | *GaussHard20d* |
|---|---|---|---|
| **RF1**$_{\text{SVM}}$ | **88%; 83%** **0.71** | **86%; 81%** **0.69** | **61%; 59%** **0.56** |
| **RF3**$_{\text{CART}}$ | 81%; 76% 0.67 | 79%; 73% 0.65 | 57%; 56% 0.54 |
| AdaBoost | 79%; 76% 0.66 | 80%; 75% 0.66 | 59%; 58% 0.55 |
| RankBoost | 84%; 80% 0.69 | 84%; 79% 0.68 | 59%; 59% 0.56 |
| RankSVM | 72%; 62% 0.53 | 70%; 55% 0.51 | 55%; 49% ; 0.56 |
| RankLS | 77%; 76% 0.67 | 77%; 73% 0.65 | 60%; 60% 0.56 |
| P − norm push | 74%; 72% 0.64 | 75%; 72% 0.64 | 60%; 57% 0.55 |
| KLR | 77%; 76% 0.67 | 79%; 75% 0.67 | 60%; 59% 0.56 |

**Table 8** Comparison of TREERANK and some competitors on simulated data. Performance is measured through the hit ratios for proportions 10% and 20%, respectively given in that order, the average precision, given below.

| $\widehat{\mathbf{Instab}}$ | *GaussEasy20d* | *GaussMed20d* | *GaussHard20d* |
|---|---|---|---|
| RF1$_{\text{SVM}}$ | *5.5* | *6.4* | *9.6* |
| RF3$_{\text{CART}}$ | *3.5* | *3.5* | *3.6* |
| AdaBoost | **2.6** | **2.9** | **3.8** |
| RankBoost | *9.2* | *11.4* | *20.8* |
| RankSVM | *4.2* | *5.0* | *12.5* |
| RankLS | *3.1* | *11.2* | *10.9* |
| P − norm push | *19.3* | *16.2* | *21.9* |
| KLR | **2.8** | *3.4* | *10.3* |

**Table 9** Comparison of TREERANK and some competitors on simulated data in terms of instability, computed with $B = 20$ bootstrap samples.

| AUC / AUC@20% / AUC@10% / AUC@5% | $\text{RF1}_{\text{SVM}}$ | $\text{RF3}_{\text{CART}}$ | AdaBoost | RankBoost | RankSVM | RankLS | P − norm push | KLR |
|---|---|---|---|---|---|---|---|---|
| *Congressional Vote* | **0.987 (±0.024)** | **0.992 (±0.015)** | 0.935 (±0.048) | 0.937 (±0.049) | 0.935 (±0.048) | 0.934 (±0.048) | 0.947 (±0.017) | 0.935 (±0.048) |
|  | **0.477 (±0.054)** | **0.495 (±0.059)** | 0.332 (±0.025) | 0.333 (±0.025) | 0.333 (±0.025) | 0.332 (±0.025) | 0.337 (±0.011) | 0.332 (±0.025) |
|  | 0.201 (±0.042) | 0.201 (±0.042) | 0.168 (±0.012) | 0.169 (±0.012) | 0.17 (±0.013) | 0.168 (±0.012) | 0.171 (±0.006) | 0.169 (±0.012) |
|  | 0.094 (±0) | 0.094 (±0) | 0.084 (±0.006) | 0.084 (±0.006) | 0.084 (±0.006) | 0.084 (±0.007) | 0.085 (±0.004) | 0.084 (±0.006) |
| *Australian Credit* | 0.926 (±0.031) | 0.934 (±0.031) | 0.928 (±0.024) | **0.937 (±0.023)** | 0.92 (±0.028) | 0.929 (±0.036) | 0.924 (±0.043) | 0.929 (±0.036) |
|  | **0.425 (±0.012)** | **0.429 (±0.014)** | 0.411 (±0.017) | 0.412 (±0.014) | 0.404 (±0.024) | 0.405 (±0.024) | 0.409 (±0.015) | 0.411 (±0.022) |
|  | 0.248 (±0.039) | 0.273 (±0.059) | 0.201 (±0.013) | 0.206 (±0.013) | 0.204 (±0.013) | 0.199 (±0.014) | 0.206 (±0.012) | 0.204 (±0.015) |
|  | 0.111 (±0.002) | 0.122 (±0.031) | 0.101 (±0.006) | 0.103 (±0.011) | 0.103 (±0.010) | 0.096 (±0.013) | 0.102 (±0.008) | 0.102 (±0.012) |
| *German Credit* | **0.793 (±0.022)** | 0.769 (±0.029) | 0.784 (±0.032) | 0.781 (±0.030) | 0.737 (±0.036) | 0.775 (±0.024) | 0.761 (±0.024) | 0.775 (±0.022) |
|  | **0.266 (±0.012)** | 0.260 (±0.020) | 0.250 (±0.021) | 0.254 (±0.016) | 0.254 (±0.012) | 0.257 (±0.013) | 0.245 (±0.015) | 0.254 (±0.016) |
|  | 0.148 (±0.020) | 0.144 (±0.022) | 0.127 (±0.013) | 0.130 (±0.010) | 0.130 (±0.004) | 0.13 (±0.008) | 0.125 (±0.010) | 0.130 (±0.011) |
|  | 0.081 (±0.021) | 0.083 (±0.029) | 0.063 (±0.009) | 0.067 (±0.005) | 0.065 (±0.005) | 0.065 (±0.003) | 0.062 (±0.006) | 0.066 (±0.005) |
| *Japanese Credit* | 0.923 (±0.038) | **0.931 (±0.042)** | 0.922 (±0.038) | **0.930 (±0.040)** | 0.904 (±0.046) | 0.910 (±0.047) | 0.897 (±0.040) | 0.916 (±0.039) |
|  | 0.414 (±0.020) | **0.417 (±0.023)** | 0.395 (±0.021) | 0.395 (±0.025) | 0.38 (±0.043) | 0.383 (±0.053) | 0.383 (±0.035) | 0.393 (±0.032) |
|  | 0.256 (±0.068) | 0.260 (±0.073) | 0.198 (±0.015) | 0.199 (±0.017) | 0.188 (±0.023) | 0.188 (±0.026) | 0.200 (±0.016) | 0.197 (±0.020) |
|  | 0.125 (±0.032) | 0.141 (±0.047) | 0.098 (±0.014) | 0.095 (±0.010) | 0.094 (±0.017) | 0.091 (±0.018) | 0.101 (±0.009) | 0.099 (±0.015) |
| *Autos MPG* | **0.978 (±0.016)** | 0.972 (±0.016) | 0.955 (±0.014) | 0.955 (±0.012) | 0.95 (±0.014) | 0.942 (±0.025) | 0.951 (±0.014) | 0.951 (±0.016) |
|  | **0.433 (±0.058)** | 0.421 (±0.016) | 0.354 (±0.005) | 0.354 (±0.005) | 0.353 (±0.008) | 0.350 (±0.011) | 0.355 (±0.008) | 0.354 (±0.008) |
|  | 0.238 (±0.079) | 0.240 (±0.08) | 0.180 (±0.006) | 0.178 (±0.006) | 0.177 (±0.005) | 0.177 (±0.006) | 0.175 (±0.004) | 0.175 (±0.005) |
|  | 0.095 (±0) | 0.095 (±0) | 0.09 (±0.005) | 0.09 (±0.006) | 0.09 (±0.003) | 0.089 (±0.006) | 0.089 (±0.003) | 0.091 (±0.004) |
| *Ionosphere* | **0.990 (±0.013)** | 0.968 (±0.032) | 0.930 (±0.015) | 0.933 (±0.013) | 0.872 (±0.081) | 0.875 (±0.055) | 0.754 (±0.048) | 0.872 (±0.062) |
|  | **0.494 (±0.062)** | 0.392 (±0.086) | 0.285 (±0.008) | 0.288 (±0.005) | 0.263 (±0.044) | 0.263 (±0.029) | 0.224 (±0.039) | 0.256 (±0.040) |
|  | 0.156 (±0) | 0.178 (±0.044) | 0.143 (±0.004) | 0.144 (±0.003) | 0.131 (±0.024) | 0.126 (±0.023) | 0.105 (±0.026) | 0.123 (±0.032) |
|  | 0.078 (±0) | 0.078 (±0) | 0.073 (±0.003) | 0.072 (±0.003) | 0.065 (±0.014) | 0.062 (±0.016) | 0.050 (±0.015) | 0.060 (±0.018) |
| *Breast Cancer Diagnosis* | **0.994 (±0.008)** | **0.988 (±0.007)** | 0.983 (±0.003) | 0.982 (±0.004) | 0.98 (±0.005) | 0.982 (±0.008) | 0.969 (±0.016) | 0.983 (±0.005) |
|  | **0.595 (±0.047)** | **0.528 (±0.005)** | 0.505 (±0.013) | 0.507 (±0.012) | 0.502 (±0.013) | 0.506 (±0.013) | 0.501 (±0.013) | 0.50 (±0.009) |
|  | 0.317 (±0.103) | 0.453 (±0.036) | 0.255 (±0.010) | 0.253 (±0.007) | 0.253 (±0.006) | 0.253 (±0.012) | 0.252 (±0.011) | 0.259 (±0.010) |
|  | 0.134 (±0) | 0.0135 (±0) | 0.130 (±0.009) | 0.125 (±0.008) | 0.126 (±0.010) | 0.130 (±0.010) | 0.128 (±0.008) | 0.131 (±0.011) |
| *Breast Cancer Original* | **0.995 (±0.006)** | 0.993 (±0.006) | 0.981 (±0.007) | 0.983 (±0.006) | 0.983 (±0.003) | 0.984 (±0.004) | 0.984 (±0.005) | 0.984 (±0.004) |
|  | **0.559 (±0.010)** | 0.556 (±0.012) | 0.536 (±0.022) | 0.534 (±0.018) | 0.537 (±0.017) | 0.537 (±0.010) | 0.531 (±0.011) | 0.539 (±0.011) |
|  | 0.442 (±0.076) | 0.407 (±0.078) | 0.259 (±0.013) | 0.265 (±0.012) | 0.271 (±0.009) | 0.265 (±0.008) | 0.268 (±0.010) | 0.264 (±0.008) |
|  | 0.146 (±0.010) | 0.142 (±0.002) | 0.131 (±0.011) | 0.132 (±0.014) | 0.137 (±0.012) | 0.134 (±0.011) | 0.136 (±0.012) | 0.136 (±0.014) |
| *Heart Disease* | **0.917 (±0.032)** | 0.911 (±0.034) | 0.879 (±0.041) | 0.885 (±0.037) | 0.883 (±0.034) | 0.891 (±0.032) | 0.876 (±0.040) | 0.892 (±0.044) |
|  | **0.416 (±0.027)** | 0.415 (±0.026) | 0.358 (±0.037) | 0.361 (±0.041) | 0.371 (±0.035) | 0.371 (±0.033) | 0.366 (±0.039) | 0.377 (±0.025) |
|  | 0.273 (±0.070) | 0.272 (±0.070) | 0.180 (±0.024) | 0.176 (±0.027) | 0.188 (±0.022) | 0.187 (±0.020) | 0.183 (±0.024) | 0.189 (±0.017) |
|  | 0.118 (±0.017) | 0.121 (±0.031) | 0.090 (±0.015) | 0.089 (±0.017) | 0.094 (±0.011) | 0.094 (±0.01) | 0.092 (±0.012) | 0.096 (±0.010) |
| *Hepatitis* | 0.864 (±0.139) | **0.868 (±0.162)** | 0.831 (±0.138) | 0.798 (±0.168) | 0.812 (±0.158) | 0.803 (±0.166) | 0.828 (±0.153) | 0.826 (±0.152) |
|  | 0.572 (±0.240) | **0.625 (±0.230)** | 0.544 (±0.205) | 0.504 (±0.225) | 0.526 (±0.248) | 0.543 (±0.231) | 0.561 (±0.234) | 0.543 (±0.230) |
|  | 0.413 (±0.138) | 0.346 (±0.162) | 0.284 (±0.103) | 0.263 (±0.115) | 0.272 (±0.125) | 0.280 (±0.116) | 0.289 (±0.118) | 0.281 (±0.116) |
|  | 0.269 (±0.190) | 0.342 (±0.154) | 0.143 (±0.051) | 0.133 (±0.057) | 0.137 (±0.062) | 0.141 (±0.057) | 0.145 (±0.058) | 0.141 (±0.057) |

**Table 10** Comparison of TREERANK and several competitors on benchmark data. Total AUC and Local AUC for three different values $u = 20, 10, 5\%$ are given in that order. Standard type error $\hat{\sigma}^2$ for these values is indicated in parentheses.

| HR@10% ; HR@20%<br>AP | RF1$_{\text{SVM}}$ | RF3$_{\text{CART}}$ | AdaBoost | RankBoost | RankSVM | RankLS | P − norm push | KLR |
|---|---|---|---|---|---|---|---|---|
| *Congressional Vote* | 100%; 100%<br>0.84 | **100%; 100%**<br>**0.85** | **100%; 100%**<br>**0.85** | 100%; 100%<br>0.84 | 100%; 100%<br>0.84 | 100%; 100%<br>0.84 | 100%; 100%<br>0.84 | **100%**; ; 100%<br>**0.85** |
| *Australian Credit* | 100%; 100%<br>0.75 | 100%; 100%<br>0.76 | 100%; 100%<br>0.75 | **100%; 100%**<br>**0.78** | 100%; 100%<br>0.76 | 100%; 100%<br>0.76 | 100%; 100%<br>0.76 | 100%; 100%<br>0.77 |
| *German Credit* | 98%; 94%<br>0.86 | **99%; 92%**<br>0.86 | 95%; 92%<br>0.86 | 96%; 93%<br>**0.87** | 97%; 94%<br>0.84 | **98%; 94%**<br>0.86 | 93%; 92%<br>0.85 | 96%; 93%<br>0.84 |
| *Japanese Credit* | 95%; 95%<br>**0.71** | 97%; 95%<br>**0.71** | 95%; 95%<br>0.71 | **97%; 96%**<br>**0.71** | 97%; 91%<br>0.69 | 94%; 92%<br>0.70 | 95%; 93%<br>0.70 | 95%; 95%<br>**0.71** |
| *Autos MPG* | **100%; 100%**<br>**0.80** | 100%; 100%<br>0.79 | 100%; 100%<br>0.77 | 100%; 100%<br>0.77 | 100%; 100%<br>0.77 | 100%; 99%<br>0.76 | 100%; 100%<br>0.77 | 100%; 100%<br>0.78 |
| *Ionosphere* | **100%; 100%**<br>**0.94** | 100%; 99%<br>0.92 | 100%; 99%<br>0.92 | 100%; 100%<br>0.92 | 94%; 93%<br>0.87 | 89%; 94%<br>0.87 | 79%; 84%<br>0.78 | 94%; 92%<br>0.86 |
| *Breast Cancer Diagnosis* | 100%; 100%<br>0.70 | 100%; 100%<br>0.70 | 100%; 100%<br>0.69 | 100%; 100%<br>0.69 | 100%; 99%<br>0.70 | 100%; 100%<br>0.69 | 100%; 100%<br>0.68 | **100%; 100%**<br>**0.71** |
| *Breast Cancer Original* | **100%; 99%**<br>**0.65** | 99%; 99%<br>;0.64 | 97%; 98%<br>0.64 | 100%; 99%<br>0.63 | 100%; 99%<br>0.65 | 100%; 99%<br>0.65 | 100%; 98%<br>0.65 | **100%; 99%**<br>**0.65** |
| *Heart Disease* | **97%; 94%**<br>0.73 | 93%; 94%<br>**0.75** | 93%; 92%<br>0.73 | 93%; 90%<br>0.71 | 97%; 94%<br>0.73 | 97%; 92%<br>0.74 | 93%; 92%<br>0.85 | 97%; 92%<br>**0.75** |
| *Hepatitis* | 65%; 40%<br>0.33 | 85%; 42%<br>0.43 | 45%; 38%<br>**0.48** | 45%; 42%<br>0.33 | 65%; 37%<br>0.33 | **85%; 47%**<br>0.42 | 65%; 42%<br>0.37 | 65%; 42%<br>0.35 |

**Table 11** Comparison of TREERANK and several competitors on benchmark data.Performance is provided here in terms of hit ratios for proportions 10% and 20% respectively (figures on top) and also in terms of average precision (figure on bottom).

26

## References

1. S. Clémençon, N. Vayatis, (2009), Tree-based Ranking Methods, IEEE Transactions on Information Theory, 9:4316-4336.
2. S. Clémençon, M. Depecker, N. Vayatis, (2011), Adaptive Partitioning Schemes for Bipartite Ranking, Machine Learning Journal, 43(1):3169.
3. S. Clémençon, M. Depecker, N. Vayatis, (2009), Bagging Ranking Trees, In Proceedings of ICMLA.
4. S. Clémençon, N. Vayatis, (2010), Ranking Forests, To be published.
5. Y. Freund, R.Iyer, R.E. Schapire, Y. Singer, (2003), An Efficient Boosting Algorithm for Combining Preferences, Journal of Machine Learning Research, 4:933-969.
6. T. Hastie, R. Tibshirani, (1990), Generalized Additive Models, Chapman & Hall/CRC.
7. J. Zhu, T. Hastie, (2005), Kernel Logistic Regression and the Import Vector Machine, Journal of Computational and Graphical Statistics, 14(1):185-205.
8. J. Friedman, T. Hastie, R. Tibshirani, (2000), Additive Logistic Regression: a Statistical View of Boosting, Annals of Statistics, 28(2):337-407.
9. T. Joachims, (2002), Optimizing Search Engines using Clickthrough Data, Proceedings of the eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 133-142.
10. T. Pahikkala, E. Tsivtsivadze, A. Airola, J. Boberg, T. Salakoski, (2007), Learning to Rank with Pairwise Regularized Least-Squares, Proceedings of SIGIR 2007 Workshop on Learning to Rank for Information Retrieval, 27-33.
11. C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, G. Hullender, (2005), Learning to Rank using Gradient Descent, Proceedings of the 22nd International Conference on Machine Learning, 89-96.
12. L. Dodd, M. Pepe, (2003), Partial AUC Estimation and Regression, Biometrics, 59(3).
13. S. Clémençon and N. Vayatis, (2007), Ranking the Best Instances, Journal of Machine Learning Research, 8:2671-2699.
14. S. Clémençon, N. Vayatis, (2008), Empirical performance maximization for linear rank statistics, Proceedings of NIPS'08, 305-312.
15. C. Rudin, (2009), The P-Norm Push: A Simple Convex Ranking Algorithm that Concentrates at the Top of the List, Journal of Machine Learning Research, 10:2233-2271.
16. S. Robertson, H. Zaragoza, (2007), On rank-based effectiveness measures and optimization, Information Retrieval, 10(3):321-339.
17. P. Bartlett, M. Jordan, J. McAuliffe, (2006), Convexity, classification and risk bounds, Journal of the American Statistical Association, 101(473):138-156.
18. P. Bartlett, A. Tewari, (2007), Sparseness vs Estimating Conditional Probabilities: Some Asymptotic Results, Journal of Machine Learning Research, 8:775-790.
19. D. Mease, A. Wyner, (2008), Evidence Contrary to the Statistical View of Boosting, Journal of Machine Learning Research, 9:131-156.
20. , Devroye, L., L. Györfi, G. Lugosi, (1996), A Probabilistic Theory of Pattern Recognition, Springer-Verlag.
21. S. Clémençon, N. Vayatis, (2010), Overlaying classifiers: a practical approach for optimal scoring, Constructive Approximation, 32(3):619-648.
22. S. Boucheron, O. Bousquet, G. Lugosi, (2005), Theory of classification: a survey of recent advances, ESAIM: Probability and Statistics, 9:323-375.
23. J. Hanley, J. McNeil, (1982), The Meaning and Use of the Area Under a ROC curve, Radiology, 143:29-36.
24. S. Clémençon, G. Lugosi, N. Vayatis, (2008), Ranking and Empirical Risk Minimization of U-statistics, The Annals of Statistics, 36:844-874.
25. N. Ailon and M. Mohri, (2010), Preference-based learning to rank, Machine Learning Journal, 80, 2:pp. 189–211.
26. L. Breiman, J. Friedman, R. Olshen, C. Stone, (1984), Classification, Regression Trees, Wadsworth and Brooks.
27. F.R. Bach, D.Heckerman, Eric Horvitz, (2006), Considering cost asymmetry in learning classifiers, Journal of Machine Learning Research, 7:1713-1741.