# A New Instance-Based Label Ranking Approach Using the Mallows Model

Weiwei Cheng and Eyke Hüllermeier

Mathematics and Computer Science University of Marburg, Germany {cheng,eyke}@mathematik.uni-marburg.de

**Abstract.** In this paper, we introduce a new instance-based approach to the label ranking problem. This approach is based on a probability model on rankings which is known as the Mallows model in statistics. Probabilistic modeling provides the basis for a theoretically sound prediction procedure in the form of maximum likelihood estimation. Moreover, it allows for complementing predictions by diverse types of statistical information, for example regarding the reliability of an estimation. Empirical experiments show that our approach is competitive to start-of-the-art methods for label ranking and performs quite well even in the case of incomplete ranking information.

 ${\bf Key \ words: \ Instance-based \ learning, \ Label \ ranking, \ Classification, \ Maximum \ likelihood \ estimation }$ 

### 1 Introduction

The topic of learning preferences has attracted increasing attention in the recent machine learning literature [1]. Label ranking, a particular preference learning scenario, studies the problem of learning a mapping from instances to rankings over a finite number of predefined labels. It can be considered as a natural generalization of the conventional classification problem, where only a single label is requested instead of a ranking of all labels.

Various approaches for label ranking have been proposed in recent years. Typically, these are extensions of learning algorithms used in binary classification problems. Ranking by pairwise comparison (RPC) is a natural extension of pairwise classification, in which binary preference models are learned for each pair of labels, and the predictions of these models are combined into a ranking of all labels [1]. Two other approaches, constraint classification (CC) and log-linear models for label ranking (LL), seek to learn linear utility functions for each individual label instead of preference predicates for pairs of labels [2,3].

In this paper, we are interested in an alternative to model-based approaches, namely the use of an *instance-based* approach. Instance-based or case-based learning algorithms have been applied successfully in various fields, such as machine learning and pattern recognition, for a long time [4]. These algorithms simply store the training data, or at least a selection thereof, and defer the processing of this data until an estimation for a new instance is requested, a property distinguishing them from typical model-based approaches. Instancebased approaches therefore have a number of potential advantages, especially in the context of the label ranking problem.

As a particular advantage of delayed processing, these learning methods may estimate the target function *locally* instead of inducing a global prediction model for the entire input domain (instance space) X. Predictions are typically obtained using only a small, locally restricted subset of the entire training data, namely those examples that are close to the query  $x \in X$  (hence X must be endowed with a distance measure). These examples are then *aggregated* in a reasonable way. As aggregating a finite set of objects from an output space  $\Omega$  is often much simpler than representing a complete  $X \to \Omega$  mapping in an explicit way, instance-based methods are especially appealing if  $\Omega$  has a complex structure.

In label ranking,  $\Omega$  corresponds to the set of all rankings of an underlying label set  $\mathcal{L}$ . To represent an  $\Omega$ -valued mapping, the aforementioned model-based approaches encode this mapping in terms of conventional binary models, either by a large set of such models in the original label space  $\mathcal{L}$  (RPC), or by a single binary model in an expanded, high-dimensional space (CC, LL). Since for instance-based methods, there is no need to represent an  $\mathbb{X} \to \Omega$  mapping explicitly, such methods can operate on the original target space  $\Omega$  directly.

The paper is organized as follows: In Section 2, we introduce the problem of label ranking in a more formal way. The core idea of our instance-based approach to label ranking, namely maximum likelihood estimation based on a special probability model for rankings, is discussed in Section 4. The model itself is introduced beforehand in Section 3. Section 5 is devoted to experimental results. The paper ends with concluding remarks in Section 6.

## 2 Label Ranking

Label ranking can be seen as an extension of the conventional setting of classification. Roughly speaking, the former is obtained from the latter through replacing single class labels by complete label rankings. So, instead of associating every instance  $\boldsymbol{x}$  from an instance space  $\mathbb{X}$  with one among a finite set of class labels  $\mathcal{L} = \{\lambda_1 \dots \lambda_n\}$ , we now associate  $\boldsymbol{x}$  with a total order of the class labels, that is, a complete, transitive, and asymmetric relation  $\succ_{\boldsymbol{x}}$  on  $\mathcal{L}$  where  $\lambda_i \succ_{\boldsymbol{x}} \lambda_j$ indicates that  $\lambda_i$  precedes  $\lambda_j$  in the ranking associated with  $\boldsymbol{x}$ . It follows that a ranking can be considered as a special type of preference relation, and therefore we shall also say that  $\lambda_i \succ_{\boldsymbol{x}} \lambda_j$  indicates that  $\lambda_i$  is *preferred* to  $\lambda_j$  given the instance  $\boldsymbol{x}$ . To illustrate, suppose that instances are students (characterized by attributes such as sex, age, and major subjects in secondary school) and  $\succ$  is a preference relation on a fixed set of study fields such as Math, CS, Physics.

Formally, a ranking  $\succ_{\boldsymbol{x}}$  can be identified with a permutation  $\pi_{\boldsymbol{x}}$  of the set  $\{1 \dots n\}$ . It is convenient to define  $\pi_{\boldsymbol{x}}$  such that  $\pi_{\boldsymbol{x}}(i) = \pi_{\boldsymbol{x}}(\lambda_i)$  is the position of  $\lambda_i$  in the ranking. This permutation encodes the (ground truth) ranking:

$$\lambda_{\pi_{\boldsymbol{x}}^{-1}(1)} \succ_{\boldsymbol{x}} \lambda_{\pi_{\boldsymbol{x}}^{-1}(2)} \succ_{\boldsymbol{x}} \ldots \succ_{\boldsymbol{x}} \lambda_{\pi_{\boldsymbol{x}}^{-1}(n)} ,$$

where  $\pi_x^{-1}(j)$  is the index of the label at position j in the ranking. The class of permutations of  $\{1 \dots n\}$  (the symmetric group of order n) is denoted by  $\Omega$ . By abuse of terminology, though justified in light of the above one-to-one correspondence, we refer to elements  $\pi \in \Omega$  as both permutations and rankings.

In analogy with the classification setting, we do not assume that there exists a deterministic  $\mathbb{X} \to \Omega$  mapping. Instead, every instance is associated with a *probability distribution* over  $\Omega$ . This means that, for each  $x \in \mathbb{X}$ , there exists a probability distribution  $\Pr(\cdot | x)$  such that, for every  $\pi \in \Omega$ ,

$$\Pr(\pi \,|\, \boldsymbol{x}) \tag{1}$$

is the probability that  $\pi_{\boldsymbol{x}} = \pi$ .

The goal in label ranking is to learn a "label ranker" in the form of an  $\mathbb{X} \to \Omega$  mapping. As training data, a label ranker uses a set of instances  $x_k, k = 1 \dots m$ , together with information about the associated rankings  $\pi_{x_k}$ . Ideally, complete rankings are given as training information. From a practical point of view, however, it is also important to allow for incomplete information in the form of a ranking

$$\lambda_{\pi_{\boldsymbol{x}}^{-1}(i_1)} \succ_{\boldsymbol{x}} \lambda_{\pi_{\boldsymbol{x}}^{-1}(i_2)} \succ_{\boldsymbol{x}} \dots \succ_{\boldsymbol{x}} \lambda_{\pi_{\boldsymbol{x}}^{-1}(i_k)}$$

where  $\{i_1, i_2 \ldots i_k\}$  is a subset of the index set  $\{1 \ldots n\}$  such that  $1 \le i_1 < i_2 < \ldots < i_k \le n$ . For example, for an instance  $\boldsymbol{x}$ , it might be known that  $\lambda_2 \succ_{\boldsymbol{x}} \lambda_1 \succ_{\boldsymbol{x}} \lambda_5$ , while no preference information is given about the labels  $\lambda_3$  or  $\lambda_4$ .

To evaluate the predictive performance of a label ranker, a suitable loss function on  $\Omega$  is needed. In the statistical literature, several distance measures for rankings have been proposed. One commonly used measure is the number of discordant pairs,

$$D(\pi, \sigma) = \{ (i, j) | i < j, \pi(i) > \pi(j) \text{ and } \sigma(i) < \sigma(j) \} , \qquad (2)$$

which is closely related to the Kendall's tau coefficient. In fact, the latter is a normalization of (2) to the interval [-1,1] that can be interpreted as a correlation measure (it assumes the value 1 if  $\sigma = \pi$  and the value -1 if  $\sigma$  is the reversal of  $\pi$ ). Kendall's tau is a natural, intuitive, and easily interpretable measure [5]. We shall focus on (2) throughout the paper, even though other distance measures could of course be used. A desirable property of any distance  $D(\cdot)$  is its invariance toward a renumbering of the elements (renaming of labels). This property is equivalent to the *right invariance* of  $D(\cdot)$ , namely  $D(\sigma\nu, \pi\nu) = D(\sigma, \pi)$  for all  $\sigma, \pi, \nu \in \Omega$ , where  $\sigma\nu = \sigma \circ \nu$  denotes the permutation  $i \mapsto \sigma(\nu(i))$ . The distance (2) is right-invariant, and so are most other commonly used metrics on  $\Omega$ .

## 3 The Mallows Model

So far, we did not make any assumptions about the probability measure (1) despite its existence. To become more concrete, we resort to a distance-based

probability model introduced by Mallows [5]. The standard Mallows model is a two-parameter model that belongs to the exponential family:

$$\Pr(\sigma \mid \theta, \pi) = \frac{\exp(\theta D(\pi, \sigma))}{\phi(\theta, \pi)},$$
(3)

where the two parameters are the location parameter (modal ranking, center ranking)  $\pi \in \Omega$  and the spread parameter  $\theta \leq 0$ . For right-invariant metrics, it can be shown that the normalization constant does not depend on  $\pi$  and, therefore, can be written as a function  $\phi(\theta)$  of  $\theta$  alone. This is due to

$$\begin{split} \phi(\theta,\pi) &= \sum_{\sigma \in \Omega} \exp(\theta D(\sigma,\pi)) = \sum_{\sigma \in \Omega} \exp(\theta D(\sigma\pi^{-1},e)) \\ &= \sum_{\sigma' \in \Omega} \exp(\theta D(\sigma',e)) = \phi(\theta) \end{split}$$

where  $e = (1 \dots n)$  is the identity ranking. More specifically, it can be shown that the normalization constant is given by [6]

$$\phi(\theta) = \prod_{j=1}^{n} \frac{1 - \exp(j\theta)}{1 - \exp(\theta)},\tag{4}$$

and that the expected distance from the center is

$$\mathbb{E}\left[D(\sigma,\pi) \,|\, \theta,\pi\right] = \frac{n \exp(\theta)}{1 - \exp(\theta)} - \sum_{j=1}^{n} \frac{j \exp(j\theta)}{1 - \exp(j\theta)} \ . \tag{5}$$

Obviously, the Mallows model assigns the maximum probability to the center ranking  $\pi$ . The larger the distance  $D(\sigma, \pi)$ , the smaller the probability of  $\sigma$ becomes. The spread parameter  $\theta$  determines how quickly the probability decreases, i.e., how peaked the distribution is around  $\pi$ . For  $\theta = 0$ , the uniform distribution is obtained, while for  $\theta \to -\infty$ , the distribution converges to the one-point distribution that assigns probability 1 to  $\pi$  and 0 to all other rankings.

## 4 Learning and Inference

Coming back to the label ranking problem and the idea of instance-based learning, consider a query instance  $\boldsymbol{x} \in \mathbb{X}$  and let  $\boldsymbol{x}_1 \dots \boldsymbol{x}_k$  denote the nearest neighbors of  $\boldsymbol{x}$  (according to an underlying distance measure on  $\mathbb{X}$ ) in the training set, where  $k \in \mathbb{N}$  is a fixed integer. Moreover, let  $\sigma_1 \dots \sigma_k \in \Omega$  denote the rankings associated, respectively, with  $\boldsymbol{x}_1 \dots \boldsymbol{x}_k$ .

In analogy to the conventional settings of classification and regression, in which the nearest neighbor estimation principle has been applied for a long time, we assume that the probability distribution  $Pr(\cdot | \boldsymbol{x})$  on  $\Omega$  is (at least approximately) *locally constant* around the query  $\boldsymbol{x}$ . By furthermore assuming

independence of the observations, the probability to observe  $\boldsymbol{\sigma} = \{\sigma_1 \dots \sigma_k\}$  given the parameters  $(\theta, \pi)$  becomes

$$\Pr(\boldsymbol{\sigma} \mid \boldsymbol{\theta}, \pi) = \prod_{i=1}^{k} \Pr(\sigma_i \mid \boldsymbol{\theta}, \pi) = \prod_{i=1}^{k} \frac{\exp\left(\boldsymbol{\theta} D\left(\sigma_i, \pi\right)\right)}{\phi(\boldsymbol{\theta})}$$
$$= \frac{\exp\left(\boldsymbol{\theta} \sum_{i=1}^{k} D\left(\sigma_i, \pi\right)\right)}{\left(\prod_{j=1}^{n} \frac{1-\exp\left(j\boldsymbol{\theta}\right)}{1-\exp\left(\boldsymbol{\theta}\right)}\right)^k}.$$
(6)

The maximum likelihood estimation (MLE) of  $(\theta, \pi)$  is then given by those parameters that maximize this probability. It is easily verified that the MLE of  $\pi$  is given by

$$\hat{\pi} = \arg\min_{\pi} \sum_{i=1}^{k} D(\sigma_i, \pi), \tag{7}$$

i.e., by the (generalized) median of the rankings  $\sigma_1 \ldots \sigma_k$ . Moreover, the MLE of  $\theta$  is derived from the average observed distance from  $\hat{\pi}$ , which is an estimation of the expected distance  $\mathbb{E}[D(\sigma, \pi)|\theta, \pi]$ :

$$\frac{1}{k}\sum_{i=1}^{k}D(\sigma_i,\hat{\pi}) = \frac{n\exp(\theta)}{1-\exp(\theta)} - \sum_{j=1}^{n}\frac{j\exp(j\theta)}{1-\exp(j\theta)}.$$
(8)

Since the right-hand side of (8) is monotone increasing, a standard line search quickly converges to the MLE [6].

Now, consider the more general case of incomplete preference information, which means that a ranking  $\sigma_i$  does not necessarily contain all labels. The probability of  $\sigma_i$  is then given by

$$\Pr(E(\sigma_i)) = \sum_{\sigma \in E(\sigma_i)} \Pr(\sigma \mid \theta, \pi) ,$$

where  $E(\sigma_i)$  denotes the set of all *consistent extensions* of  $\sigma_i$ : A permutation  $\sigma \in \Omega$  is a consistent extension of  $\sigma$  if it ranks all labels that also occur in  $\sigma_i$  in the same order.

The probability of observing the neighbor rankings  $\boldsymbol{\sigma} = (\sigma_1 \dots \sigma_k)$  then becomes

$$\Pr(\boldsymbol{\sigma} \mid \boldsymbol{\theta}, \pi) = \prod_{i=1}^{k} \Pr(E(\sigma_i) \mid \boldsymbol{\theta}, \pi) = \prod_{i=1}^{k} \sum_{\sigma \in E(\sigma_i)} \Pr(\sigma \mid \boldsymbol{\theta}, \pi)$$
$$= \frac{\prod_{i=1}^{k} \sum_{\sigma \in E(\sigma_i)} \exp\left(\boldsymbol{\theta} D(\sigma, \pi)\right)}{\left(\prod_{j=1}^{n} \frac{1 - \exp(j\boldsymbol{\theta})}{1 - \exp(\boldsymbol{\theta})}\right)^k} .$$
(9)

Computing the MLE of  $(\theta, \pi)$  by maximizing this probability now becomes more difficult. For label sets of small to moderate size, say up to 7, one can afford a

simple brute force approach, namely an exhaustive search over  $\Omega$  to find the center ranking  $\pi$ , combined with a numerical procedure to optimize the spread  $\theta$ . For larger label sets, this procedure becomes too inefficient. Here, we propose an approximation algorithm that can be seen as an instance of the EM (Expectation-Maximization) family of algorithms.

The algorithm works as follows. Starting from an initial (complete) center ranking  $\hat{\pi}$ , each incomplete neighbor ranking  $\sigma_i$  is replaced by the most probable consistent extension given  $\hat{\pi}$ . Regardless of  $\theta$ , this extension is obviously given by a ranking in  $\arg\min_{\sigma\in E(\sigma_i)} D(\sigma, \hat{\pi})$ . It can be found by (minimally) re-ranking the center  $\hat{\pi}$  so as to make it consistent with the incomplete ranking  $\sigma_i$ . Having replaced all neighbor rankings by their most probable extensions, an MLE ( $\theta, \pi$ ) can be derived as described for the case of complete information above. The center ranking  $\hat{\pi}$  is then replaced by  $\pi$ , and the whole procedure is iterated until the center does not change any more. In the following, we discuss two subproblems of the algorithm in more detail, namely the solution of the median problem (7), which needs to be solved to find an MLE  $\pi$ , and the choice of an initial center ranking.

Solving the (generalized) median problem (7) is known to be NP-complete for Kendall's tau, i.e., if the distance D is given by the number of rank inversions [7]. To solve this problem approximately, we make use of the fact that Kendall's tau is well approximated by Spearman's rank correlation [8], and that the median can be computed for this measure (i.e., for D given by the sum of squared rank differences) by a procedure called *Borda count* [9]: Given a (complete) ranking  $\sigma_i$  of n labels, the top-label receives n votes, the second-ranked n-1 votes, and so on. Given k rankings  $\sigma_1 \ldots \sigma_k$ , the sum of the k votes are computed for each label, and the labels are then ranked according to their total votes.

The choice of the initial center ranking in the above algorithm is of course critical. To find a good initialization, we again resort to the idea of solving the problem (7) approximately using the Borda count principle. At the beginning, however, the neighbor rankings  $\sigma_k$  are still incomplete. To handle this situation, we make the simplifying assumption that the completions are uniformly distributed in  $E(\sigma_i)$ . Again, this is an approximation, since we actually proceed from the Mallows and not from the uniform model. On the basis of this assumption, we can show the following result (proof omitted due to space restrictions).

**Theorem 1.** Let a set of incomplete rankings  $\sigma_1 \ldots \sigma_k$  be given, and suppose the associated complete rankings  $S_1 \ldots S_k$  to be distributed, respectively, uniformly in  $E(\sigma_1) \ldots E(\sigma_k)$ . The expected sum of distances  $D(\pi, S_1) + \ldots + D(\pi, S_k)$ , with D the sum of squared rank distances, becomes minimal for the ranking  $\pi$  which is obtained by a generalized Borda count, namely a Borda count with a generalized distribution of votes from incomplete rankings: If  $\sigma_i$  is an incomplete ranking of  $m \leq n$  labels, then the label on rank  $i \in \{1 \ldots m\}$  receives (m-i+1)(n+1)/(m+1) votes, while each missing label receives a vote of (n + 1)/2.

Table 1. Statistics for the semi-synthetic and real datasets

dataset	#examples	#classes	#features
iris	150	3	4
$_{\rm wine}$	178	3	13
$_{\rm glass}$	214	6	9
vehicle	846	4	18
dtt	2465	4	24
$\operatorname{cold}$	2465	4	24

## 5 Experimental Results

#### 5.1 Methods

In this section, we compare our instance-based (nearest neighbor, NN) approach with existing methods for label ranking, namely ranking by pairwise comparison (RPC), constraint classification (CC), and log-linear models for label ranking (LL). Since space restrictions prevent from a detailed review, we refer to the original literature and [1] for a short review of these methods. Regarding the concrete implementation and parameterization of these methods, we also follow [1].

To fit the Mallows model, we test the two previously discussed variants, namely the exhaustive search which guarantees an optimal solution (NNE) and the approximation algorithm outlined in Section 4 (NNH). The parameter k (neighborhood size) was selected through cross validation on the training set. As a distance measure on the instance space we used the Euclidean distance (after normalizing the attributes).

#### 5.2 Data

We used two real-world data sets, dtt and cold, from the bioinformatics field. These data sets contain two types of genetic data, namely phylogenetic profiles and DNA microarray expression data for the Yeast genome.<sup>1</sup> The genome consists of 2465 genes, and each gene is represented by an associated phylogenetic profile of length 24. Using these profiles as input features, we investigated the task of predicting a "qualitative" representation of an expression profile; see [1] for a detailed description and motivation of this task.

In addition to the real-world data sets, the following multiclass datasets from the UCI repository of machine learning databases and the Statlog collection were included in the experimental evaluation: iris, wine, glass, vehicle. For each of these datasets, a corresponding ranking dataset was generated in the following manner: We trained a naive Bayes classifier on the respective dataset. Then, for each example, *all* the labels present in the dataset were ordered with respect to decreasing predicted class probabilities (in the case of ties, labels with lower

<sup>&</sup>lt;sup>1</sup> This data is publicly available at http://www1.cs.columbia.edu/compbio/.

**Table 2.** Experimental results in terms of Kendall's tau (mean and standard deviation) for different missing label rates (parameter p).

iris	0%	10%	20%	30%	40%	50%	60%	70%
RPC	$.885 \pm .068$	$.888 \pm .064$	$.886 \pm .060$	$.871 \pm .074$	$.854 \pm .082$	$.837 \pm .089$	$.779 \pm .110$	$.674 \pm .139$
CC	$.836 \pm .089$	$.825 \pm .095$	$.815 \pm .088$	$.807 \pm .099$	$.788 \pm .105$	$.766 \pm .115$	$.743 \pm .131$	.708±.105
LL	$.818 \pm .088$	$.811 \pm .089$	$.805 \pm .087$	$.806 \pm .087$	$.800 \pm .091$	$.788 \pm .087$	$.778 \pm .096$	$.739 \pm .186$
NNE	$.960 \pm .036$	$.956 {\pm} .041$	$941 \pm .044$	$.934 {\pm} .049$	$.915 {\pm} .056$	$.882 {\pm} .085$	$.859 {\pm} .082$	$.812 \pm .107$
NNH	$.966 {\pm} .034$	$.948 \pm .036$	$.917 \pm .051$	$.863 \pm .072$	$.822 \pm .088$	$.802 \pm .084$	$.767 \pm .122$	.733±.104
wine								
RPC	$.921 \pm .053$	$.900 \pm .067$	$.886 \pm .073$	$.902 \pm .063$	$.910 \pm .065$	$.882 \pm .082$	$.864 \pm .097$	.822±.118
CC	$.933 \pm .043$	$.918 \pm .057$	$.929 \pm .058$	$.911 \pm .059$	$.922 {\pm} .057$	$.885 \pm .074$	$.853 \pm .078$	$.802 \pm .123$
LL	$.942 \pm .043$	$.944 \pm .046$	$.939 \pm .051$	$.944 {\pm} .042$	$.933 \pm .062$	$.918 \pm .065$	$.906 \pm .072$	$.864 {\pm} .094$
NNE	$.952 \pm .048$	$.945 \pm .051$	$.943 \pm .055$	$.940 \pm .054$	$.941 {\pm} .050$	$.930 {\pm} .058$	$.910 \pm .061$	$.677 \pm .173$
NNH	$.953 {\pm} .042$	$.949 {\pm} .041$	$.949 \pm .041$	$.933 \pm .048$	$.899 \pm .075$	$.709 \pm .186$	$.591 \pm .210$	$.587 \pm .180$
glass								
RPC	$.882 {\pm} .042$	$.875 {\pm} .046$	$.867 {\pm} .044$	$.851 {\pm} .052$	$.840 {\pm} .053$	$.813 {\pm} .062$	$.799 {\pm} .054$	$.754 \pm .076$
CC	$.846 \pm .045$	$.848 \pm .053$	$.838 \pm .059$	$.835 \pm .054$	$.833 \pm .051$	$.807 \pm .066$	$.789 \pm .052$	$.747 \pm .061$
LL	$.817 \pm .060$	$.815 \pm .061$	$.813 \pm .063$	$.819 \pm .062$	$.819 \pm .060$	$.809 \pm .066$	$.806 {\pm} .065$	$.807 {\pm} .063$
NNE	$.875 \pm .063$	$.866 \pm .059$	$.840 \pm .059$	$.803 \pm .062$	$.750 \pm .071$	$.677 \pm .066$	$.598 \pm .082$	$1.500 \pm .078$
NNH	$.865 \pm .059$	$.847 \pm .062$	$.810 \pm .056$	$.754 \pm .069$	$.691 \pm .063$	$.633 \pm .061$	$.550 \pm .069$	$.484 \pm .079$
vehicle								
RPC	$.854 \pm .025$	$.848 \pm .025$	$.847 \pm .024$	$.834 \pm .026$	$.823 \pm .032$	$.803 \pm .033$	$.786 \pm .036$	$.752 \pm .041$
CC	$.855 \pm .022$	$.848 \pm .026$	$.849 \pm .026$	$.839 {\pm} .025$	$.834 {\pm} .026$	$.827 {\pm} .026$	$.810 {\pm} .026$	$.791 {\pm} .030$
LL	$.770 \pm .037$	$.769 \pm .035$	$.769 \pm .033$	$.766 \pm .040$	$.770 \pm .038$	$.764 \pm .031$	$.757 \pm .038$	$.756 \pm .036$
NNE	$.863 {\pm} .030$	$.859 {\pm} .031$	$.847 \pm .029$	$.834 \pm .031$	$.822 \pm .030$	$.795 \pm .033$	$.766 \pm .034$	$1.723 \pm .036$
NNH	$.862 \pm .025$	$.852 \pm .024$	$.845 \pm .030$	$.828 \pm .029$	$.798 \pm .031$	$.776 \pm .033$	$.748 \pm .032$	$.701 \pm .047$
dtt								
RPC	$.174 \pm .034$	$.172 \pm .034$	$.168 \pm .036$	$.166 \pm .036$	$.164 \pm .034$	$.153 \pm .035$	$.144 \pm .028$	$.125 \pm .030$
CC	$.180 \pm .037$	$.178 \pm .034$	176±.033	$ .172 \pm .032 $	$.165 \pm .033$	$.158 \pm .033$	$.149 \pm .031$	$1.136 \pm .033$
LL	$.167 \pm .034$	$.168 \pm .033$	$168 \pm .034$	$.168 \pm .034$	$.167 {\pm} .033$	$.167 {\pm} .036$	$.162 \pm .032$	$.156 {\pm} .034$
NNE	$.182 \pm .036$	$.179 \pm .036$	$.173 \pm .036$	$.169 \pm .036$	$.162 \pm .036$	$.161 \pm .037$	$.154 \pm .036$	$.136 \pm .035$
NNH	$.191 {\pm} .034$	$1.183 {\pm} .037$	$ .176 \pm .036 $	$.168 \pm .038$	$.163 \pm .034$	$.146 \pm .036$	$.145 \pm .033$	$1.128 \pm .035$
cold								
RPC	$.221 {\pm} .028$	$.217 \pm .028$	$.213 \pm .030$	$.212 \pm .030$	$.208 \pm .030$	$.201 \pm .030$	$.188 \pm .030$	$.174 \pm .031$
CC	$.220 \pm .029$	$.219 \pm .030$	$.212 \pm .030$	$.212 \pm .028$	$.205 \pm .024$	$.197 \pm .030$	$.185 \pm .031$	$.162 \pm .035$
LL	$.209 \pm .028$	$.210 \pm .031$	$.206 \pm .030$	$.210 \pm .030$	$.203 \pm .031$	$.203 \pm .031$	$.202 {\pm} .032$	$.192 \pm .031$
NNE	$.230 \pm .028$	$.226 \pm .029$	$.220 \pm .030$	$.213 \pm .031$	$.199 {\pm} .029$	$.195 \pm .033$	$.190 \pm .035$	$.188 \pm .035$
NNH	$.244 {\pm} .026$	$.237 {\pm} .028$	$.235 {\pm} .031$	$.226 {\pm} .024$	$.220 {\pm} .029$	$.214 {\pm} .029$	$.199 \pm .030$	$.192 \pm .032$

index are ranked first). Thus, by substituting the single labels contained in the original multiclass datasets with the complete rankings, we obtain the label ranking datasets required for our experiments. A summary of the data sets and their properties is given in Table 1.

#### 5.3 Experiments and Results

Results were derived in terms of the Kendall's tau correlation coefficient from five repetitions of a ten-fold cross-validation. To model incomplete preferences, we modified the training data as follows: A biased coin was flipped for every label in a ranking in order to decide whether to keep or delete that label; the probability for a deletion is specified by a parameter p.

The results are summarized in Table 2. As can be seen, NN is quite competitive to the model-based approaches and often outperforms these methods. In any case, it is always close to the best result. It is also remarkable that NN seems to be quite robust toward missing preferences and compares comparably well in this regard. This was not necessarily expected, since NN uses only local information, in contrast to the other approaches that induce global models. Our approximation algorithm NNH gives very good approximations of NNE throughout and is especially appealing for large label sets: It dramatically reduces the runtime (not shown due to space restrictions) without any significant decrease of the performance.

A nice feature of our approach, not shared by the model-based methods, is that it comes with a natural measure of the reliability of a prediction. In fact, the smaller the parameter  $\theta$ , the more peaked the distribution around the center ranking and, therefore, the more reliable this ranking becomes as a prediction. To test whether (the estimation of)  $\theta$  is indeed a good measure of uncertainty of a prediction, we used it to compute a kind of *accuracy-rejection* curve: By averaging over five 10-fold cross validations (with NNE), we computed an accuracy degree  $\tau_x$  (the average Kendall's tau) and a reliability degree  $\theta_x$  for each instance x. The instances are then sorted in decreasing order of reliability. Our curve plots a value p against the mean  $\tau$ -value of the first p percent of the instances. Given that  $\theta$  is indeed a good indicator of reliability, this curve should be decreasing, because the higher p, the more instances with a less strong  $\theta$ -value are taken into consideration. As can be seen in Fig. 1, the curves obtained for our data sets are indeed decreasing and thus provide evidence for our claim that  $\theta$  may serve as a reasonable indicator of the reliability of a prediction.

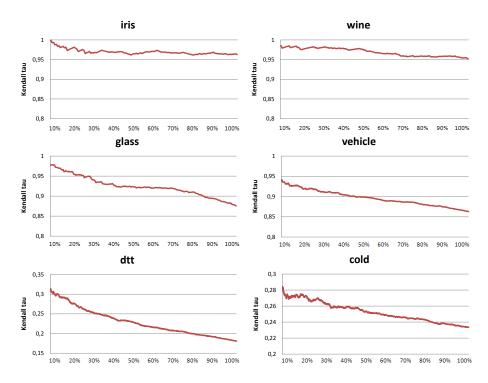


Fig. 1. Accuracy-rejection curves computed on the basis of the parameter  $\theta$ .

## 6 Conclusions and Future Work

In this paper, we have introduced an instance-based (nearest neighbor) approach to the label ranking problem that has recently attracted attention in the field of machine learning. Our basic inference principle is a consistent extension of the nearest neighbor estimation principle, as used previously for well-known learning problems such as classification and regression: Assuming that the conditional (probability) distribution of the output given the query is locally constant, we derive a maximum likelihood estimation based on the Mallows model, a special type of probability model for rankings. Our first empirical results are quite promising and suggest that this approach is fully competitive, in terms of predictive accuracy, to (model-based) state-of-the-art methods for label ranking. Besides, it has some further advantages, as it does not only produce a single ranking as an estimation but instead delivers a probability distribution over all rankings. This distribution can be used, for example, to quantify the reliability of the predicted ranking.

Currently, we are working on extensions and variants of the label ranking problem, such as calibrated label ranking and multi-label classification [10]. In fact, we believe that the approach proposed in this paper can be extended to a solid framework that not only allows for solving the label ranking problem itself but also variants thereof.

## References

- Hüllermeier, E., Fürnkranz, J., Cheng, W., Brinker, K.: Label Ranking by Learning Pairwise Preferences. Artificial Intelligence 172(16-17), 1897–1916 (2008)
- Har-Peled, S., Roth, D., Zimak, D.: Constraint Classification for Multiclass Classification and Ranking. In Becker, S., Thrun, S., Obermayer, K., eds.: Advances in Neural Information Processing Systems 15, 785-792 (2003)
- Dekel, O., Manning, C., Singer, Y.: Log-Linear Models for Label Ranking. In Touretzky, D.S., Thrun, S., Saul, L.K., Schölkopf, B., eds.: Advances in Neural Information Processing Systems 16, 497-504 (2004)
- Aha, D., Kibler, D., Albert, M.: Instance-Based Learning Algorithms. Machine Learning 6(1), 37-66 (1991)
- 5. Mallows, C.: Non-Null Ranking Models. Biometrika 44(1), 114–130 (1957)
- Fligner, M., Verducci, J.: Distance Based Ranking Models. Journal of the Royal Statistical Society 48(3), 359-369 (1986)
- Alon, N.: Ranking Tournaments. SIAM Journal on Discrete Mathematics 20(1), 134-142 (2006)
- Diaconis, P., Graham, R.: Spearman's Footrule as a Measure of Disarray. Journal of the Royal Statistical Society 39(2), 262–268 (1977)
- 9. Saari, D.: Chaotic Elections!: A Mathematician Looks at Voting. American Mathematical Society (2001)
- Brinker, K., Hüllermeier, E.: Case-based Multilabel Ranking. In: Proc. IJCAI-07, 20th International Joint Conference on Artificial Intelligence, 701-707, Hyderabad, India (January 2007)