

# Pairwise Preference Learning and Ranking

Johannes Fürnkranz<sup>1</sup> and Eyke Hüllermeier<sup>2</sup>

<sup>1</sup> Austrian Research Institute for Artificial Intelligence  
Schottengasse 3, A-1010 Wien, Austria  
juffi@oefai.at

<sup>2</sup> Informatics Institute, Marburg University  
Hans-Meerwein-Str., Lahnberge, D-35032 Marburg, Germany  
eyke@mathematik.uni-marburg.de

**Abstract.** We consider supervised learning of a ranking function, which is a mapping from instances to total orders over a set of labels (options). The training information consists of examples with partial (and possibly inconsistent) information about their associated rankings. From these, we induce a ranking function by reducing the original problem to a number of binary classification problems, one for each pair of labels. The main objective of this work is to investigate the trade-off between the quality of the induced ranking function and the computational complexity of the algorithm, both depending on the amount of preference information given for each example. To this end, we present theoretical results on the complexity of pairwise preference learning, and experimentally investigate the predictive performance of our method for different types of preference information, such as top-ranked labels and complete rankings. The domain of this study is the prediction of a rational agent's ranking of actions in an uncertain environment.

## 1 Introduction

The problem of learning with or from preferences has recently received a lot of attention within the machine learning literature.<sup>3</sup> The problem is particularly challenging because it involves the prediction of complex structures, such as weak or partial order relations, rather than single values. Moreover, training input will not, as it is usually the case, be offered in the form of complete examples but may comprise more general types of information, such as relative preferences or different kinds of indirect feedback.

More specifically, the learning scenario that we will consider in this paper consists of a collection of training examples which are associated with a finite set of decision alternatives. Following the common notation of supervised learning, we shall refer to the latter as *labels*. However, contrary to standard classification, a training example is not assigned a single label, but a set of *pairwise preferences* between labels, expressing that one label is preferred over another.

The goal is to use these pairwise preferences for predicting a total order, a *ranking*, of all possible labels for a new training example. More generally, we seek to induce a

---

<sup>3</sup> Space restrictions prevent a thorough review of related work in this paper, but we refer the reader to (Fürnkranz and Hüllermeier, 2003).

*ranking function* that maps instances (examples) to rankings over a fixed set of decision alternatives (labels), in analogy to a *classification function* that maps instances to single labels. To this end, we investigate the use of *round robin learning* or *pairwise classification*. As will be seen, round robin appears particularly appealing in this context since it can be extended from classification to preference learning in a quite natural manner.

The paper is organized as follows: In the next section, we introduce the learning problem in a formal way. The extension of pairwise classification to pairwise preference learning and its application to ranking are discussed in section 3. Section 4 provides some results on the computational complexity of pairwise preference learning. Results of several experimental studies investigating the predictive performance of our approach under various training conditions are presented in section 5. We conclude the paper with some final remarks in section 6.

## 2 Learning Problem

We consider the following learning problem:

---

**Given:**

- a set of *labels*  $L = \{\lambda_i \mid i = 1 \dots c\}$
- a set of *examples*  $E = \{e_k \mid k = 1 \dots n\}$
- for each training example  $e_k$ :
  - a set of *preferences*  $P_k \subseteq L \times L$ , where  $(\lambda_i, \lambda_j) \in P_k$  indicates that label  $\lambda_i$  is preferred over label  $\lambda_j$  for example  $e_k$  (written as  $\lambda_i \succ_k \lambda_j$ )

**Find:** a function that orders the labels  $\lambda_i, i = 1 \dots c$  for any given example.

---

This setting has been previously introduced as *constraint classification* by Harpeled et al. (2002). As has been pointed out in their work, the above framework is a generalization of several common learning settings, in particular (see *ibidem* for a formal derivation of these and other results)

- *ranking*: Each training example is associated with a total order of the labels, i.e., for each pair of labels  $(\lambda_i, \lambda_j)$  either  $\lambda_i \succ_k \lambda_j$  or  $\lambda_j \succ_k \lambda_i$  holds.
- *classification*: A single class label  $\lambda_i$  is assigned to each example. This implicitly defines the set of preferences  $\{\lambda_i \succ_k \lambda_j \mid 1 \leq j \neq i \leq c\}$ .
- *multi-label classification*: Each training example  $e_k$  is associated with a subset  $S_k \subseteq L$  of possible labels. This implicitly defines the set of preferences  $\{\lambda_i \succ_k \lambda_j \mid \lambda_i \in S_k, \lambda_j \in L \setminus S_k\}$ .

As pointed out before, we will be interested in predicting a ranking (total order) of the labels. Thus, we assume that for each instance, there exists a total order of the labels, i.e., the pairwise preferences form a transitive and asymmetric relation. For many practical applications, this assumption appears to be acceptable at least for the *true* preferences. Still, more often than not the observed or *revealed* preferences will be incomplete or inconsistent. Therefore, we do not require the *data* to be consistent in the sense that transitivity and asymmetry applies to the  $P_k$ . We only assume that  $P_k$  is irreflexive ( $\lambda_i \not\succeq \lambda_i$ ) and anti-symmetric ( $\lambda_i \succ \lambda_j \Rightarrow \lambda_j \not\succeq \lambda_i$ ). (Note that  $0 \leq |P_k| \leq c(c-1)/2$  as a consequence of the last two properties.)

### 3 Pairwise Preference Ranking

A key idea of our approach is to learn a separate theory for each of the  $c(c-1)/2$  pairwise preferences between two labels. More formally, for each possible pair of labels  $(\lambda_i, \lambda_j)$ ,  $1 \leq i < j \leq c$ , we learn a model  $M_{ij}$  that decides for any given example whether  $\lambda_i \succ \lambda_j$  or  $\lambda_j \succ \lambda_i$  holds. The model is trained with all examples  $e_k$  for which either  $\lambda_i \succ_k \lambda_j$  or  $\lambda_j \succ_k \lambda_i$  is known. All examples for which nothing is known about the preference between  $\lambda_i$  and  $\lambda_j$  are ignored.

At classification time, an example is submitted to all  $c(c-1)/2$  theories. If classifier  $M_{ij}$  predicts  $\lambda_i \succ \lambda_j$ , we count this as a vote for  $\lambda_i$ . Conversely, the prediction  $\lambda_j \succ \lambda_i$  would be considered as a vote for  $\lambda_j$ . The labels are ranked according to the number of votes they receive from all models  $M_{ij}$ . Ties are first broken according to the frequency of the labels in the top rank (the class distribution in the classification setting) and then randomly.

We refer to the above technique as *pairwise preference ranking* or *round robin ranking*. It is a straight-forward generalization of pairwise or one-against-one classification, aka round robin learning, which solves multi-class problems by learning a separate theory for each pair of classes. In previous work, Fürnkranz (2002) showed that, for rule learning algorithms, this technique is preferable to the more commonly used one-against-all classification method, which learns one theory for each class, using the examples of this class as positive examples and all others as negative examples. Interestingly, despite its complexity being quadratic in the number of classes, the algorithm is no slower than the conventional one-against-all technique (Fürnkranz, 2002). We will generalize these results in the next section.

### 4 Complexity

Consider a learning problem with  $n$  training examples and  $c$  labels.

**Theorem 1.** *The total number of training examples over all  $c(c-1)/2$  binary preference learning problems is*

$$\sum_{k=1}^n |P_k| \leq n \max_k |P_k| \leq n \binom{c}{2} = n \frac{c(c-1)}{2}$$

*Proof.* Each of the  $n$  training examples  $e_k$  will be added to all  $|P_k|$  binary training sets that correspond to one of its preferences  $\lambda_i \succ_k \lambda_j$ . Thus, the total number of training examples is  $\sum_{k=1}^n |P_k|$ . As the number of preferences for each example is bounded from above by  $\max_k |P_k|$ , this number is no larger than  $n \max_k |P_k|$ , which in turn is bounded from above by the size of a complete set of preferences  $nc(c-1)/2$ .  $\square$

**Corollary 1.** (Fürnkranz, 2002) *For a classification problem, the total number of training examples is only linear in the number of classes.*

*Proof.* A class label expands to  $c-1$  preferences, therefore  $\sum_{k=1}^n |P_k| = (c-1)n$ .  $\square$

Note that we only considered the number of training examples, but not the complexity of the learner that runs on these examples. For an algorithm with a linear run-time complexity  $O(n)$  it follows immediately that the total run-time is  $O(dn)$ , where  $d$  is the maximum (or average) number of preferences given for each training example. For a learner with a super-linear complexity  $O(n^a)$ ,  $a > 1$ , the total run-time is much lower than  $O((dn)^a)$  because the training effort is not spent on one large training set, but on many small training sets. In particular, for a complete preference set, the total complexity is  $O(c^2n^a)$ , whereas the complexity for  $d = c - 1$  (round robin classification) is only  $O(cn^a)$  (Fürnkranz, 2002).

For comparison, the only other technique for learning in this setting that we know of (Har-Peled et al., 2002) constructs twice as many training examples (one positive and one negative for each preference of each example), and these examples are projected into a space that has  $c$  times as many attributes as the original space. Moreover, all examples are put into a single training set for which a separating hyper-plane has to be found. Thus, under the (reasonable) assumption that an increase in the number of features has approximately the same effect as a corresponding increase in the number of examples, the total complexity becomes  $O((cdn)^a)$  if the algorithm for finding the separating hyper-plane has complexity  $O(n^a)$  for a two-class training set of size  $n$ .

In summary, the overall complexity of pairwise constraint classification depends on the number of known preferences for each training example. While being quadratic in the number of labels if a complete ranking is given, it is only linear for the classification setting. In any case, it is more efficient than the technique proposed by Har-Peled et al. (2002). However, it should be noted that the price to pay is the large number of classifiers that have to be stored and tested at classification time.

## 5 Empirical Results

The previous sections have shown that round robin learning can be extended to induce a ranking function from a set of preferences instead of a single label. Yet, it turned out that computational complexity might become an issue. Especially, since a ranking induces a quadratic number of pairwise preferences, the complexity for round robin ranking becomes quadratic in the number of labels. In this context, one might ask whether it could be possible to improve efficiency at the cost of a tolerable decrease in performance: Could the learning process perhaps ignore some of the preferences without decreasing predictive accuracy too much? Apart from that, incomplete training data is clearly a point of practical relevance, since complete rankings will rarely be observable.

The experimental evaluation presented in this section is meant to investigate issues related to incomplete training data in more detail, especially to increase our understanding about the trade-off between the number of pairwise preferences available in the training data and the quality of the learned ranking function. For a systematic investigation of questions of such kind, we need data for which, in principle, a complete ranking is known for each example. This information allows a systematic variation of the amount of preference information in the training data, and a precise evaluation of the predicted rankings on the test data. Since we are not aware of any suitable real-world datasets, we have conducted our experiments with synthetic data.

## 5.1 Synthetic Data

We consider the problem of learning the ranking function of an expected utility maximizing agent. More specifically, we proceed from a standard setting of expected utility theory:  $A = \{a_1, \dots, a_c\}$  is a set of actions the agent can choose from and  $\Omega = \{\omega_1, \dots, \omega_m\}$  is a set of world states. The agent faces a problem of *decision under risk* where decision consequences are lotteries: Choosing act  $a_i$  in state  $\omega_j$  yields a utility of  $u_{ij} \in \mathbb{R}$ , where the probability of state  $\omega_j$  is  $p_j$ . Thus, the *expected utility* of act  $a_i$  is given by

$$\mathbb{E}(a_i) = \sum_{j=1}^m p_j \cdot u_{ij}. \quad (1)$$

Expected utility theory justifies (1) as a criterion for ranking actions and, hence, gives rise to the following preference relation:

$$a_i \succ a_j \Leftrightarrow \mathbb{E}(a_i) > \mathbb{E}(a_j). \quad (2)$$

Now, suppose the probability vector  $p = (p_1, \dots, p_m)$  to be a parameter of the decision problem (while  $A, \Omega$  and the utility matrix  $U = (u_{ij})$  are fixed).

The above decision-theoretic setting can be used for generating synthetic data for preference learning. The set of instances corresponds to the set of probability vectors  $p$ , which are generated at random according to a uniform distribution over  $\{p \in \mathbb{R}^m \mid p \geq 0, p_1 + \dots + p_m = 1\}$ . The ranking function associated with an example is given by the ranking defined in (2). Thus, an experiment is characterized by the following parameters: The number of actions/labels ( $c$ ), the number of world states ( $m$ ), the number of examples ( $n$ ), and the utility matrix which is generated at random through independent and uniformly distributed entries  $u_{ij} \in [0, 1]$ .

## 5.2 Experimental Setup

In the following, we will report on results of experiments with ten different states ( $m = 10$ ) and various numbers of labels ( $c = 5, 10, 20$ ). For each of the three configurations we generated ten different data sets, each one originating from a different randomly chosen utility matrix  $U$ . The data sets consisted of 1000 training and 1000 test examples. For each example, the data sets provided the probability vector  $p \in \mathbb{R}^m$  and a complete ranking of the  $c$  possible actions.<sup>4</sup> The training examples were labeled with a subset of the complete set of pairwise preferences as imposed by the ranking in the data set. The subsets that were selected for the experiments are described one by one for the experiments.

We used the decision tree learner C4.5 (Quinlan, 1993) in its default settings<sup>5</sup> to learn a model for each pairwise preference. For all examples in the test set we obtained a final ranking using simple voting and tie breaking as described in section 3.

<sup>4</sup> The occurrence of actions with equal expected utility has probability 0.

<sup>5</sup> Our choice of C4.5 as the learner was solely based on its versatility and wide availability. If we had aimed at maximizing performance on this particular problem, we would have chosen an algorithm that can directly represent the separating hyperplanes for each binary preference.

The predicted ranks were then compared with the actual ranks. Our primary evaluation measures were the error rate of the top rank (for comparing classifications) and the Spearman rank correlation coefficient (for comparing complete rankings).

### 5.3 Ranking vs. Classification

Figure 1 shows experimental results for (a) using the full set of  $c(c - 1)/2$  pairwise preferences, (b) for the classification setting which uses only the  $c - 1$  preferences that involve the top label, and (c) for the complementary setting with the  $(c - 1)(c - 2)/2$  preferences that do *not* involve the top label. There are several interesting things to note for these results. First, the difference between the error rates of the classification and the ranking setting is comparably small. Thus, if we are only interested in the top rank, it may often suffice to use the pairwise preferences that involve the top label. The advantage in this case is of course the reduced complexity which becomes linear in the number of labels. On the other hand, the results also show that the complete ranking information can be used to improve classification accuracy, at least if this information is available for each training example and if one is willing to pay the price of a quadratic complexity.

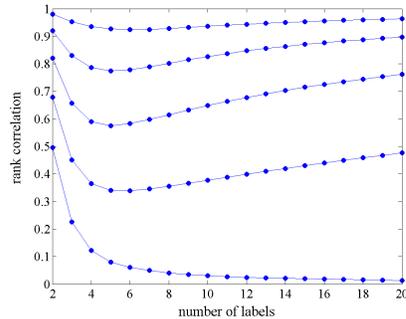
The results for the complementary setting show that the information of the top rank preferences is crucial: When dropping this information and using only those pairwise preferences that do not involve the top label, the error rate on the top rank increases considerably, and is much higher than the error rate for the classification setting. This is a bit surprising if we consider that in the classification setting, the average number of training examples for learning a model  $M_{ij}$  is much smaller than in the complementary setting. Interestingly, the effective number of training examples for the top labels might nevertheless decrease. In fact, in our learning scenario we will often have a few *dominating* actions whose utility degrees are systematically larger than those of other actions. In the worst case, the same action is optimal for all probability vectors  $p$ , and the complementary set will not contain any information about it. While this situation is of course rather extreme, the class distribution is indeed very unbalanced in our scenario. For example, we determined experimentally for  $c = m = 10$  and  $n = 1000$  that the probability of having the same optimal action for more than half of the examples is  $\approx 2/3$ , and that the expected Gini-index of the class distribution is  $\approx 1/2$ .

With respect to the prediction of complete rankings, the performance for learning from the complementary set of preferences is almost as good as the performance for learning from the complete set of preferences, whereas the performance of the ranking induced from the classification setting is considerably worse. This time, however, the result is hardly surprising and can easily be explained by the amount of information provided in the two cases. In fact, the complementary set determines the ranking of  $c - 1$  among the  $c$  labels, whereas the top label alone does hardly provide any information about the complete ranking.

As another interesting finding note that the classification accuracy decreases with an increasing number of labels, whereas the rank correlation increases (this is also revealed by the curves in Figure 3 below). In other words, the quality of the predicted rankings increases, even though the quality of the predictions for the individual ranks decreases. This effect can first of all be explained by the fact that the (classification)

$c$	prefs	error	rank corr.
5	ranking	$13.380 \pm 8.016$	$0.907 \pm 0.038$
	classification	$14.400 \pm 8.262$	$0.783 \pm 0.145$
	complement	$32.650 \pm 14.615$	$0.872 \pm 0.051$
10	ranking	$15.820 \pm 8.506$	$0.940 \pm 0.018$
	classification	$16.670 \pm 9.549$	$0.711 \pm 0.108$
	complement	$24.310 \pm 9.995$	$0.937 \pm 0.018$
20	ranking	$24.030 \pm 4.251$	$0.966 \pm 0.004$
	classification	$26.370 \pm 5.147$	$0.697 \pm 0.066$
	complement	$32.300 \pm 3.264$	$0.966 \pm 0.004$

**Fig. 1.** Comparison of ranking (a complete set of preferences is given) vs. classification (only the preferences for the top rank are given). Also shown are the results for the complementary setting (all preferences for the top rank are omitted).



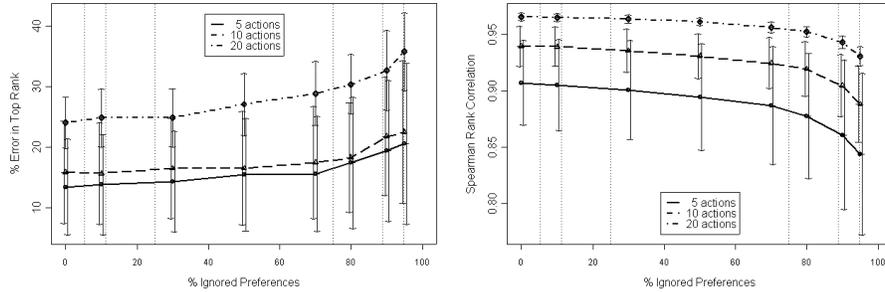
**Fig. 2.** Expected Spearman rank correlation as a function of the number of labels if all models  $M_{ij}$  have an error rate of  $\epsilon$  (curves are shown for  $\epsilon = 0.1, 0.2, 0.3, 0.4, 0.5$ ).

error is much more affected by an increase of the number of labels. As an illustration, consider random guessing: The chances of guessing the top label correctly are  $1/m$ , whereas the expected value of the rank correlation is 0 regardless of  $m$ . Moreover, one might speculate that the importance of a correct vote of each individual model  $M_{ij}$  decreases with an increasing number of labels. Roughly speaking, incorrect classifications of individual learners are better compensated on average. This conjecture is also supported by an independent experiment in which we simulated a set of homogeneous models  $M_{ij}$  through biased coin flipping with a prespecified error rate. It turned out that the quality measures for predicted rankings tend to increase if the number of labels becomes large (though the dependence of the measures on the number of labels is not necessarily monotone, see Fig. 2).

#### 5.4 Missing Preferences

While the previous results shed some light on the trade-off between utility and costs for two special types of preference information, namely top-ranked labels and complete rankings, they do not give a satisfactory answer for the general case. The selected set of preferences in the classification setting is strongly focused on a particular label for each example, thus resulting in a very biased distribution. In the following, we will look at the quality of predicted rankings when selecting random subsets of pairwise preferences from the full sets with equal right.

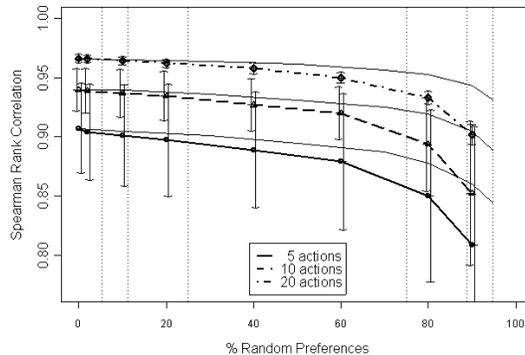
Figure 3 shows the curves for the classification error in the top rank and the average Spearman rank correlation of the predicted and the true ranking over the number of preferences. To generate these curves, we started with the full set of preferences, and ignored increasingly larger fractions of it. This was implemented with a parameter  $p_i$  that caused any given preference in the training data to be ignored with probability  $p_i$  ( $100 \times p_i$  is plotted on the  $x$ -axis).



**Fig. 3.** Average error rate (left) and Spearman rank correlation (right) for various percentages of ignored preferences. The error bars indicate the standard deviations. The vertical dotted lines on the right indicate the number of preferences for classification problems (for 5,10, and 20 classes), those on the left are the complementary sizes.

The similar shape of the three curves (for 5, 10, and 20 labels) suggests that the decrease in the ranking quality can be attributed solely to the missing preferences while it seems to be independent of the number of labels. In particular, one is inclined to conclude that—contrary to the case where we focused on the top rank—it is in general *not* possible to reduce the number of training preferences by an order of magnitude (i.e., from quadratic to linear in the number of labels) without severely decreasing the ranking quality. This can also be seen from the three dotted vertical lines in the right half of the graphs. These lines indicate the percentage of preferences that were present in the classification setting for 5, 10, and 20 labels (from inner-most to outer-most). A comparison of the error rates, given by the intersection of a line with the corresponding curve, to the respective error rates in Figure 1 shows an extreme difference between the coincidental selection of pairwise preferences and the systematic selection which is focused on the top rank.

Nevertheless, one can also see that about half of the preferences can be ignored while still maintaining a reasonable performance level. Even though it is quite common that learning curves are concave functions of the size of the training set, the descent in accuracy appears to be remarkably flat in our case. One might be tempted to attribute this to the redundancy of the pairwise preferences induced by a ranking: In principle, a ranking  $\rho$  could already be reconstructed from the  $c - 1$  preferences  $\rho_1 \succ \rho_2, \dots, \rho_{c-1} \succ \rho_c$ , which means that only a small fraction of the pairwise preferences are actually needed. Still, one should be careful with this explanation. First, we are not trying to reconstruct a single ranking but rather to solve a slightly different problem, namely to learn a ranking function. Second, our learning algorithm does actually not “reconstruct” a ranking as suggested above. In fact, our simple voting procedure does not take the dependencies between individual models  $M_{ij}$  into account, which means that these models do not really cooperate. On the contrary, what the voting procedure exploits is just the redundancy of preference information: The top rank is the winner only because it is preferred in  $c - 1$  out of the  $c(c - 1)/2$  pairwise comparisons.



**Fig. 4.** Average Spearman rank correlation over various percentages of random preferences. The error bars indicate the standard deviations. The solid thin lines are the curves for ignored preferences (Figure 3).

Finally, note that the shape of the curves probably also depends on the number of training examples. We have not yet investigated this issue because we were mainly interested in the possibility of reducing the complexity by more than a constant factor without losing too much of predictive accuracy. It would be interesting, for example, to compare (a) using  $p\%$  of the training examples with full preferences and (b) using all training examples with  $p\%$  of the pairwise preferences.

### 5.5 Mislabeled Preferences

Recall that our learning scenario assumes preference structures to be complete rankings of labels, that is transitive and asymmetric relations. As already pointed out, we do not make this assumption for *observed* preferences: First, we may not have access to complete sets of preferences (the case studied in the previous section). Second, the process generating the preferences might reproduce the underlying total order incorrectly and, hence, produce inconsistent preferences. The latter problem is quite common, for example, in the case of human judgments.

To simulate this behavior, we adopted the following model: Proceeding from the pairwise preferences induced by a given ranking, a preference  $\lambda_i \succ \lambda_j$  was kept with probability  $1 - p_s$ , whereas with probability  $p_s$ , one of the preferences  $\lambda_i \succ \lambda_j$  and  $\lambda_j \succ \lambda_i$  was selected by a coin flip. Thus, in approximately  $p_s/2$  cases, the preference will point into the wrong direction.<sup>6</sup> For  $p_s = 0$ , the data remain unchanged, whereas the preferences in the training data are completely random for  $p_s = 1$ .

Figure 4 shows the average Spearman rank correlations that were observed in this experiment. Note that the shape of the curve is almost the same as the shape of the curves for ignored preferences. It is possible to directly compare these two curves because in both graphs a level of  $n\%$  means that  $100 - n\%$  of the preferences are still

<sup>6</sup> In fact, we implemented the procedure by selecting  $p_s/2$  preferences and reversing their sign.

intact. The main difference is that in Figure 3, the remaining  $n\%$  of the preferences have been ignored, while in Figure 4 they have been re-assigned at random. To facilitate this comparison, we plotted the curves for ignored preferences (the same ones as in Figure 3) into the graph (with solid, thin lines).

It is interesting to see that in both cases the performance degrades very slowly at the beginning, albeit somewhat steeper than if the examples are completely ignored. Roughly speaking, completely omitting a pairwise preference appears to be better than including a random preference. This could reasonably be explained by the learning behavior of a classifier  $M_{ij}$ : If  $M_{ij}$  does already perform well, an additional correct example will probably be classified correctly and thus improve  $M_{ij}$  only slightly (in decision tree induction, for example,  $M_{ij}$  will even remain completely unchanged if the new example is classified correctly). As opposed to this, an incorrect example will probably be classified incorrectly and thus produce a more far-reaching modification of  $M_{ij}$  (in decision tree induction, an erroneous example might produce a completely different tree). All in all, the “expected benefit” of  $M_{ij}$  caused by a random preference is negative, whereas it is 0 if the preference is simply ignored.

From this consideration one may conclude that a pairwise preference should better be ignored if it is no more confident than a coin flip. This can also be grasped intuitively, since the preference does not provide any information in this case. If it is more confident, however, it clearly carries some information and it might then be better to include it, even though the best way of action will still depend on the number and reliability of the preferences already available. Note that our experiments do not suggest any strategy for deciding whether or not to include an *individual* preference, given information about the uncertainty of that preference. In our case, each preference is equally uncertain. Thus, the only reasonable strategies are to include all of them or to ignore the complete sample. Of course, the first strategy will be better as soon as the probability of correctness exceeds  $1/2$ , and this is also confirmed by the experimental results. For example, the correlation coefficient remains visibly above 0.8 even if 80% of the preferences are assigned by chance and, hence, the probability of a particular preference to be correct is only 0.6. One may conjecture that pairwise preference ranking is particularly robust toward noise, since an erroneous example affects only a single classifier  $M_{ij}$  which in turn has a limited influence on the eventually predicted ranking.

## 6 Concluding Remarks

We have introduced pairwise preference learning as an extension of pairwise classification to constraint classification, a learning scenario where training examples are labeled with a preference relation over all possible labels instead of a single class label as in the conventional classification setting. From this information, we also learn one model for each pair of classes, but focus on learning a complete ranking of all labels instead of only predicting the most likely label. Our main interest was to investigate the trade-off between ranking quality and the amount of training information (in terms of the number of preferences that are available for each example). We experimentally investigated this trade-off by varying parameters of a synthetic domain that simulates a decision-theoretic agent which ranks its possible actions according to an unknown

utility function. Roughly speaking, the results show that large parts of the information about pairwise preferences can be ignored in round robin ranking without losing too much predictive performance. In the classification setting, where one is only interested in predicting the top label, it also turned out that using the full ranking information rather than restricting to the pairwise preferences involving the top label does even improve the classification accuracy, suggesting that the lower ranks do contain valuable information. For reasons of efficiency, however, it might still be advisable to concentrate on the smaller set of preferences, thereby reducing the size of the training set by an order of magnitude.

The main limitation of our technique is probably the assumption of having enough training examples for learning each pairwise preference. For data with a very large number of labels and a rather small set of preferences per example, our technique will hardly be applicable. In particular, it is unlikely to be successful in collaborative filtering problems (Goldberg et al., 1992; Resnick and Varian, 1997; Breese et al., 1998), although these can be mapped onto the constraint classification framework in a straightforward way. A further limitation is the quadratic number of theories that has to be stored in memory and evaluated at classification time. However, the increase in memory requirements is balanced by an increase in computational efficiency in comparison to the technique of Har-Peled et al. (2002). In addition, pairwise preference learning inherits many advantages of pairwise classification, in particular its implementation can easily be parallelized because of its reduction to independent subproblems. Finally, we have assumed an underlying total order of the items which needs to be recovered from partial observations of preferences. However, partial orders (cases where several labels are equally preferred) may also occur in practical applications. We have not yet investigated the issue of how to generate (and evaluate) partial orders from learned pairwise predictions. Similarly, our current framework does not provide a facility for discriminating between cases where we know that a pair of labels is of equal preference and cases where we don't know anything about their relative preferences.

There are several directions for future work. First of all, it is likely that the prediction of rankings can be improved by combining the individual models' votes in a more sophisticated way. Several authors have looked at techniques for combining the predictions of pairwise theories into a final ranking of the available options. Proposals include weighting the predicted preferences with the classifiers' confidences (Fürnkranz, 2003) or using an iterative algorithm for combining pairwise probability estimates (Hastie and Tibshirani, 1998). However, none of the previous works have evaluated their techniques in a ranking context, and some more elaborate proposals, like error-correcting output decoding (Allwein et al., 2000), organizing the pairwise classifiers in a tree-like structure (Platt et al., 2000), or using a stacked classifier (Savicky and Fürnkranz, 2003) are specifically tailored to a classification setting. Taking into account the fact that we are explicitly seeking a ranking could lead to promising alternatives. For example, we are thinking about selecting the ranking which minimizes the number of predicted preferences that need to be reversed in order to make the predicted relation transitive. Departing from the counting of votes might also offer possibilities for extending our method to the prediction of preference structures more general than rankings (total orders), such as weak preference relations where some of the labels might not be comparable.

Apart from theoretical considerations, an important aspect of future work concerns the practical application of our method and its evaluation using real-world problems. Unfortunately, real-world data sets that fit our framework seem to be quite rare. In fact, currently we are not aware of any data set of significant size that provides instances in attribute-value representation plus an associated complete ranking over a limited number of labels.

### Acknowledgments

Johannes Fürnkranz is supported by an APART stipend (no. 10814) of the *Austrian Academy of Sciences*. The Austrian Research Institute for Artificial Intelligence is supported by the Austrian Federal Ministry of Education, Science and Culture.

### References

- E. L. Allwein, R. E. Schapire, and Y. Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
- J. S. Breese, D. Heckerman, and C. Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In G. F. Cooper and S. Moral (eds.), *Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence (UAI-98)*, pp. 43–52, Madison, WI, 1998. Morgan Kaufmann.
- J. Fürnkranz. Round robin classification. *Journal of Machine Learning Research*, 2: 721–747, 2002.
- J. Fürnkranz. Round robin ensembles. *Intelligent Data Analysis*, 7(5), 2003. To appear.
- J. Fürnkranz and E. Hüllermeier. Pairwise preference learning and ranking. Technical Report OEFAI-TR-2003-14, Austrian Research Institute for Artificial Intelligence, Wien, Austria, 2003.
- D. Goldberg, D. Nichols, B. M. Oki, and D. Terry. Using collaborative filtering to weave and information tapestry. *Communications of the ACM*, 35(12):61–70, 1992.
- S. Har-Peled, D. Roth, and D. Zimak. Constraint classification: A new approach to multiclass classification. In N. Cesa-Bianchi, M. Numao, and R. Reischuk (eds.), *Proceedings of the 13th International Conference on Algorithmic Learning Theory (ALT-02)*, pp. 365–379, Lübeck, Germany, 2002. Springer-Verlag.
- T. Hastie and R. Tibshirani. Classification by pairwise coupling. In M. Jordan, M. Kearns, and S. Solla (eds.), *Advances in Neural Information Processing Systems 10 (NIPS-97)*, pp. 507–513. MIT Press, 1998.
- J. C. Platt, N. Cristianini, and J. Shawe-Taylor. Large margin DAGs for multiclass classification. In S. A. Solla, T. K. Leen, and K.-R. Müller (eds.), *Advances in Neural Information Processing Systems 12 (NIPS-99)*, pp. 547–553. MIT Press, 2000.
- J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, CA, 1993.
- P. Resnick and H. R. Varian. Special issue on recommender systems. *Communications of the ACM*, 40(3), 1997.
- P. Savicky and J. Fürnkranz. Combining pairwise classifiers with stacking. In *Advances in Intelligent Data Analysis: Proceedings of the 5th International Symposium (IDA-03)*, Berlin, Germany, 2003. Springer-Verlag. To appear.