

Preference Rules for Label Ranking: Mining Patterns in Multi-Target Relations

Cláudio Rebelo de Sá^{a,b,*}, Paulo Azevedo^d, Carlos Soares^c,
Alípio Mário Jorge^e, Arno Knobbe^b

^a*INESC TEC, Porto, Portugal*

^b*LIACS, Leiden, Netherlands*

^c*INESC TEC, Faculdade de Engenharia, Universidade do Porto*

^d*HasLab, INESC TEC, Departamento de Informática, Universidade do Minho*

^e*INESC TEC, Faculdade de Ciências, Universidade do Porto*

Abstract

In this paper, we investigate two variants of association rules for preference data, Label Ranking Association Rules and Pairwise Association Rules. Label Ranking Association Rules (LRAR) are the equivalent of Class Association Rules (CAR) for the Label Ranking task. In CAR, the consequent is a single class, to which the example is expected to belong to. In LRAR, the consequent is a ranking of the labels. The generation of LRAR requires special support and confidence measures to assess the similarity of rankings. In this work, we carry out a sensitivity analysis of these similarity-based measures. We want to understand which datasets benefit more from such measures and which parameters have more influence in the accuracy of the model. Furthermore, we propose an alternative type of rules, the Pairwise Association Rules (PAR), which are defined as association rules with a set of pairwise preferences in the consequent. While PAR can be used both as descriptive and predictive models, they are essentially descriptive models. Experimental results show the potential of both approaches.

Keywords: Label Ranking, Association Rules, Pairwise Comparisons

*Corresponding author

Email addresses: `claudio.r.sa@inesctec.pt` (Cláudio Rebelo de Sá),
`pja@di.uminho.pt` (Paulo Azevedo), `csoares@fe.up.pt` (Carlos Soares), `amjorge@fc.up.pt`
(Alípio Mário Jorge), `a.j.knobbe@liacs.leidenuniv.nl` (Arno Knobbe)

1. Introduction

Label ranking is a topic in the machine learning literature [1, 2, 3] that studies the problem of learning a mapping from instances to rankings over a finite number of predefined labels. One characteristic that clearly distinguishes Label Ranking problems from classification problems is the order relation between the labels. While a classifier aims at finding the true class on a given unclassified example, the label ranker will focus on the relative preferences between a set of labels/classes. These relations represent relevant information from a decision support perspective, with possible applications in various fields such as elections, dominance of certain species over the others, user preferences, etc.

Due to its intuitive representation, Association Rules [4] have become very popular in data mining and machine learning tasks (e.g. mining rankings [5], classification [6] or even Label Ranking [7, 8]). In [7], association rules were adapted for the prediction of rankings, which are referred to as Label Ranking Association Rules (LRAR). A different approach, Rule-Based Label Ranking (RBLR) [8], adapts the Dominance-based Rough Set Approach (DRSA) [9] for predicting rankings in the Label Ranking task. Both LRAR and RBLR can be used for predictive or descriptive purposes.

LRAR are relations, like typical association rules, between an antecedent and a consequent ($A \rightarrow C$), defined by interest measures. The distinction lies in the fact that the consequent is a complete ranking. Because the degree of similarity between rankings can vary, it leads to several interesting challenges. For instance, how to treat rankings that are very similar but not exactly equal. To tackle this problem, similarity-based interest measures were defined to evaluate LRAR. Such measures can be applied to existing rule generation methods [7] (e.g. APRIORI [4]).

One important issue for the use of LRAR is the threshold that determines what should and should not be considered sufficiently similar. Here we present the results of sensitivity analysis study to show how LRAR behave in different

30 scenarios, to understand the effect of this threshold better. Whether there is a rule of thumb or this threshold is data-specific is the type of questions we investigate here. Additionally we also want to understand which parameters have more influence in the predictive accuracy of the method.

Another important issue is related to the large number of distinct rankings. 35 Despite the existence of many competitive approaches in Label Ranking, Decision trees [10, 2], k -Nearest Neighbor [11, 2] or LRAR [7], problems with a large number of distinct rankings can be hard to model. One real-world example with a relatively large number of rankings, is the sushi dataset [12]. This dataset compares demographics of 5000 Japanese citizens with their preferred 40 sushi types. With only 10 labels, it has more than 4900 distinct rankings. Even though it has been known in the preference learning community for a while, no results with high predictive accuracy have been published, to the best of our knowledge. This might be due to noise in the data or simply because of inconsistency in the ratings provided by the people interviewed [13]. Cases like 45 this have motivated the appearance of new approaches, e.g. to mine ranking data [5], where association rules are used to find patterns within rankings.

We propose a method which combines the two approaches mentioned above [7, 5], to extract interesting information from datasets even when the number of different rankings is very high. We define Pairwise Association Rules (PAR) as 50 association rules with one or more pairwise comparisons in the consequent. In this work, we present an approach to identify PAR and analyze the findings in two real world datasets.

By decomposing rankings into the unitary preference relation i.e. *pairwise comparisons*, we can look for sub-ranking patterns, which are expected to be 55 more frequent.

LRAR and PAR can be regarded as a specialization of general association rules that are obtained from data containing preferences, which we refer to as *Preference Rules*. These two approaches are complementary in the sense that they can give different insights from multi-target relations that can be found in 60 preference data [14]. We use LRAR and PAR in this work as predictive and

descriptive models, respectively.

The paper is organized as follows: Sections 2 and 3 introduce the task of association rule mining and the Label Ranking problem, respectively; Section 4 describes the Label Ranking Association Rules and Section 5 the Pairwise Association Rules proposed here; Section 6 presents the experimental setup and discusses the results; finally, Section 7 concludes this paper.

2. Association Rule Mining

An association rule (AR) is an implication: $A \rightarrow C$ where $A \cap C = \emptyset$ and $A, C \subseteq desc(\mathbb{X})$, where $desc(\mathbb{X})$ is the set of descriptors of instances in the instance space \mathbb{X} , typically pairs $\langle attribute, value \rangle$. The training data is represented as $D = \{\langle x_i \rangle\}$, $i = 1, \dots, n$, where x_i is a vector containing the values $x_i^j, j = 1, \dots, m$ of m independent variables, \mathcal{A} , describing instance i . We also denote $desc(x_i)$ as the set of descriptors of instance x_i .

2.1. Interest measures

There are many interest measures to evaluate association rules [15], but typically they are characterized by *support* and *confidence*. Here, we summarize some of the most common, assuming a rule $A \rightarrow C$ in D .

Support. Percentage of the instances in D that contain A and C :

$$sup(A \rightarrow C) = \frac{\#\{x_i | A \cup C \subseteq desc(x_i), x_i \in D\}}{n}$$

Confidence. percentage of instances that contain C from the set of instances that contain A :

$$conf(A \rightarrow C) = \frac{sup(A \rightarrow C)}{sup(A)}$$

Coverage. Proportion of examples in D that contain the antecedent of a rule: *coverage* [16]:

$$coverage(A \rightarrow C) = sup(A)$$

We say that a rule $A \rightarrow C$ covers an instance x , if $A \subseteq desc(x)$.

Lift. Measures the independence of the consequent, C , relative to the antecedent, A :

$$lift(A \rightarrow C) = \frac{sup(A \rightarrow C)}{sup(A) \cdot sup(C)}$$

Lift values vary from 0 to $+\infty$. If A is independent from C then $lift(A \rightarrow C) \sim 1$.

2.2. Methods

The original method for induction of AR is the APRIORI algorithm, proposed in 1994 [4]. APRIORI identifies all AR that have support and confidence higher than a given minimal support threshold (*minsup*) and a minimal confidence threshold (*minconf*), respectively. Thus, the model generated is a set of AR, \mathcal{R} , of the form $A \rightarrow C$, where $A, C \subseteq desc(\mathbb{X})$, and $sup(A \rightarrow C) \geq minsup$ and $conf(A \rightarrow C) \geq minconf$. For a more detailed description see [4].

Despite the usefulness and simplicity of APRIORI, it runs a time consuming candidate generation process and needs substantial time and memory space, proportional to the number of possible combinations of the descriptors. Additionally it needs multiple scans of the data and typically generates a very large number of rules. Because of this, many alternative methods were previously proposed, such as hashing [17], dynamic itemset counting [18], parallel and distributed mining [19] and mining integrated into relational database systems [20].

A major breakthrough in itemset mining has been brought by the algorithm FP-Growth (Frequent pattern growth method) [21], which starts by efficiently projecting the original data base into a compact tree data structure (FP-tree). From the FP-tree, itemset support can be calculated without revisiting the original dataset, which also avoids the generation of candidate itemsets. With respect to APRIORI there is an enormous reduction both on computational time and space necessary. FP-Growth approach is also able to efficiently find long itemsets.

2.3. Pruning

AR algorithms typically generate a large number of rules (possibly tens of thousands), some of which represent only small variations from others. This is

known as the rule explosion problem [22] which should be dealt with by pruning mechanisms. Many rules must be discarded for computational and simplicity reasons.

115 Pruning methods are usually employed to reduce the amount of rules without reducing the quality of the model. For example, an AR algorithm might find rules for which the confidence is only marginally improved by adding further conditions to their antecedent. Another example is when the consequent C of a rule $A \rightarrow C$ has the same distribution independently of the antecedent A . In
 120 these cases, we should not consider these rules as meaningful.

Improvement. A common pruning method is based on the improvement that a refined rule yields in comparison to the original one [22]. The *improvement* of a rule is defined as the smallest difference between the confidence of a rule and the confidence of all sub-rules sharing the same consequent:

$$imp(A \rightarrow C) = \min(\forall A' \subset A, conf(A \rightarrow C) - conf(A' \rightarrow C))$$

125 As an example, if one defines a minimum improvement $minImp = 1\%$, the rule $A \rightarrow C$ will be kept if $conf(A \rightarrow C) - conf(A' \rightarrow C) \geq 1\%$, for any $A' \subset A$.
 If $imp(A \rightarrow C) > 0$ we say that $A \rightarrow C$ is a productive rule.

Significant rules. Another way to prune nonproductive rules is to use statistical tests [23]. A rule is *significant* if the confidence improvement over all
 130 its generalizations is statistically significant. The rule $A \rightarrow C$ is significant if $\forall A' \rightarrow C, A' \subset A$ the difference $conf(A \rightarrow C) - conf(A' \rightarrow C)$ is statistically significant for a given significance level (α).

3. Label Ranking

In Label Ranking, given an instance x from the instance space \mathbb{X} , the goal
 135 is to predict the ranking of the labels $\mathcal{L} = \{\lambda_1, \dots, \lambda_k\}$ associated with x [24]. A ranking can be represented as a *strict total order* over \mathcal{L} , defined on the permutation space Ω .

The Label Ranking task is similar to the classification task, where instead of a class we want to predict a ranking of the labels. As in classification, we do not assume the existence of a deterministic $\mathbb{X} \rightarrow \Omega$ mapping. Instead, every instance is associated with a *probability distribution* over Ω [2]. This means that, for each $x \in \mathbb{X}$, there exists a probability distribution $\mathcal{P}(\cdot|x)$ such that, for every $\pi \in \Omega$, $\mathcal{P}(\pi|x)$ is the probability that π is the ranking associated with x . The goal in Label Ranking is to learn the mapping $\mathbb{X} \rightarrow \Omega$. The training data contains a set of instances $D = \{\langle x_i, \pi_i \rangle\}$, $i = 1, \dots, n$, where x_i is a vector containing the values $x_i^j, j = 1, \dots, m$ of m independent variables, \mathcal{A} , describing instance i and π_i is the corresponding target ranking.

Rankings can be represented with total or partial orders and vice-versa.

Total orders. A *strict total order* over \mathcal{L} is defined as a binary relation, \succ , on a set \mathcal{L} [25], which is:

1. Irreflexive: $\lambda_a \not\succeq \lambda_a$
2. Transitive: $\lambda_a \succ \lambda_b$ and $\lambda_b \succ \lambda_c$ implies $\lambda_a \succ \lambda_c$
3. Asymmetric: if $\lambda_a \succ \lambda_b$ then $\lambda_b \not\succeq \lambda_a$ ¹
4. Connected: For any λ_a, λ_b in \mathcal{L} , either $\lambda_a \succ \lambda_b$ or $\lambda_b \succ \lambda_a$

A *strict ranking* [3], a *complete ranking* [27], or simply a *ranking* can be represented by a *strict total order* over \mathcal{L} . A strict total order can also be represented as a permutation π of the set $\{1, \dots, k\}$, such that $\pi(a)$ is the position, or *rank*, of λ_a in π . For example, the *strict total order* $\lambda_3 \succ \lambda_1 \succ \lambda_2 \succ \lambda_4$ can be represented as $\pi = (2, 3, 1, 4)$.

However, in real-world ranking data, we do not always have clear and unambiguous preferences, i.e. strict total orders [28]. Hence, sometimes we have to deal with *indifference* [29] and *incomparability* [30]. For illustration purposes, let us consider the scenario of elections, where a set of n voters vote on k candidates. If a voter feels that two candidates have identical proposals, then these can be expressed as indifferent so they are assigned the same rank (i.e. a tie).

¹Asymmetry can be derived from 1. and 2. [26].

To represent ties, we need a more relaxed setting, called *non-strict total orders*, or simply *total orders*, over \mathcal{L} , by replacing the binary strict order relation, \succ , with the binary partial order relation, \succeq where the following properties hold [25]:

- 170 1. Reflexive: $\lambda_a \succeq \lambda_a$
2. Transitive: $\lambda_a \succeq \lambda_b$ and $\lambda_b \succeq \lambda_c$ implies $\lambda_a \succeq \lambda_c$
3. Antisymmetric: $\lambda_a \succeq \lambda_b$ and $\lambda_b \succeq \lambda_a$ implies $\lambda_a = \lambda_b$
4. Connected: For any λ_a, λ_b in \mathcal{L} , either $\lambda_a \succeq \lambda_b$, $\lambda_b \succeq \lambda_a$ or $\lambda_b = \lambda_a$

These non-strict total orders can represent *partial rankings* (rankings with ties)
 175 [3]. For example, the *non-strict total order* $\lambda_1 \succ \lambda_2 = \lambda_3 \succ \lambda_4$ can be represented as $\pi = (1, 2, 2, 3)$.

Additionally, real-world data may lack preference data regarding two or more labels, which is known as *incomparability*. Continuing with the elections example, the lack of information about one or two of the candidates, λ_a and λ_b , leads
 180 to incomparability, $\lambda_a \perp \lambda_b$. In other words, the voter cannot decide whether the candidates are equivalent or select one as the preferred, because he does not know the candidates. Incomparability should not be confused with intrinsic properties of the objects, as if we are comparing apples and oranges. Instead, it is like trying to compare two different types of apple without ever having tried
 185 at least one of them. In this cases, we can use *partial orders*.

Partial orders. Similar to *total orders*, there are *strict* and *non-strict partial orders*. Let us consider the *non-strict partial orders* (which can also be referred to as *partial orders*) where the binary relation, \succeq , over \mathcal{L} is [25]:

1. Reflexive: $\lambda_a \succeq \lambda_a$
- 190 2. Transitive: $\lambda_a \succeq \lambda_b$ and $\lambda_b \succeq \lambda_c$ implies $\lambda_a \succeq \lambda_c$
3. Antisymmetric: $\lambda_a \succeq \lambda_b$ and $\lambda_b \succeq \lambda_a$ implies $\lambda_a = \lambda_b$

We can represent partial orders with *subrankings* [5] or *incomplete rankings* [31]. For example, the *partial order* $\lambda_1 \succ \lambda_2 \succ \lambda_4$ can be represented as $\pi = (1, 2, 0, 3)$, where 0 represents $\lambda_1, \lambda_2, \lambda_4 \perp \lambda_3$.

Several learning algorithms were proposed for modeling Label Ranking data in recent years. These can be grouped as decomposition-based or direct. *Decomposition methods* divide the problem into several simpler problems (e.g., multiple binary problems). An example is ranking by pairwise comparisons [1] and mining rank data [5]. *Direct methods* treat the rankings as target objects without any decomposition. Examples of that include decision trees [10, 2], k -Nearest Neighbors [11, 2] and the linear utility transformation [32, 33]. This second group of algorithms can be divided into two approaches. The first one contains methods that are based on statistical distributions of rankings (e.g. [2]), such as Mallows [34], or Plackett-Luce [31]. The other group of methods are based on measures of similarity or correlation between rankings (e.g. [10, 35]).

Label Ranking-specific pre-processing methods have also been proposed, e.g. MDLP-R [36] and EDiRa [37]. Both are *direct methods* and based on measures of similarity. Considering that supervised discretization approaches usually provide better results than unsupervised methods [38], such methods can be of a great importance in the field. In particular, for AR-like algorithms, such as the ones proposed in this work, which are typically not suitable for numerical data.

Below, we briefly describe some of these Label Ranking approaches (including both direct and decomposition methods) with which we compare our method in the experimental part (Section 6).

3.1.1. Rule-Based Label Ranking

Rule-Based Label Ranking (RBLR) [8] is a rule-based approach that aims to deliver interpretable results in the form of logical rules. It is essentially a decomposition method, where the rankings are decomposed into pairwise comparisons (λ_a, λ_b) and considered as a further attribute called *relation attribute* [8]. It uses an adapted version of the Dominance-based Rough Set Approach (DRSA) for Label Ranking data to transform the features into a gain and cost criteria.

3.1.2. Instance-Based Plackett-Luce

Instance-Based Plackett-Luce (IB-PL) is a local prediction method based on the nearest neighbor estimation principle [39]. Given a new instance \hat{x} it uses the $\{\pi_1, \dots, \pi_\beta\}$ rankings of the β nearest neighbors to predict the ranking $\hat{\pi}$ associated with \hat{x} . The estimation of $\hat{\pi}$ is made using the Plackett-Luce (PL) model. For the PL model, the probability to observe a ranking π is:

$$\mathcal{P}(\pi|v) = \prod_{i=1}^k \frac{v_{\pi^{-1}(i)}}{v_{\pi^{-1}(i)} + v_{\pi^{-1}(i+1)} + \dots + v_{\pi^{-1}(k)}}$$

where $v = (v_1, \dots, v_k)$ is obtained with a Maximum Likelihood Estimation and can be seen as a vector indicating the skill, score or popularity of each object [39]. The larger the parameter v_i in comparison to the remaining parameters, the higher the probability of λ_i to be on a top rank.

3.1.3. Label Ranking by Learning Pairwise Preferences

Ranking by pairwise comparisons basically consists of reducing the problem of ranking into several classification problems. In the learning phase, the original problem is formulated as a set of pairwise preferences problems. Each problem is concerned with one pair of labels of the ranking, $(\lambda_i, \lambda_j) \in \mathcal{L}, 1 \leq i < j \leq k$. The target attribute is the relative order between them, $\lambda_i \succ \lambda_j$. Then, a separate model \mathcal{M}_{ij} is obtained for each pair of labels. Considering $\mathcal{L} = \{\lambda_1, \dots, \lambda_k\}$, there will be $h = \frac{k(k-1)}{2}$ classification problems to model.

In the prediction phase, each model is applied to every pair of labels to obtain a prediction of their relative order. The predictions are then combined to derive rankings, which can be done in several ways. The simplest is to order the labels, for each example, considering the predictions of the models \mathcal{M}_{ij} as votes. This topic has been well studied and documented [40, 24].

More detailed information on Label Ranking methods can be found in [41].

3.2. Evaluation

Given an instance x_i with real ranking π_i and a ranking $\hat{\pi}_i$ predicted by a
 250 Label Ranking model, several loss functions on Ω can be used to evaluate the
 accuracy of the prediction. One such function is the number of discordant label
 pairs:

$$\mathcal{D}(\pi, \hat{\pi}) = \#\{(a, b) | \pi(a) > \pi(b) \wedge \hat{\pi}(a) < \hat{\pi}(b)\}$$

If there are no discordant label pairs, the distance $\mathcal{D} = 0$. Alternatively, the
 function to define the number of concordant pairs is:

$$\mathcal{C}(\pi, \hat{\pi}) = \#\{(a, b) | \pi(a) > \pi(b) \wedge \hat{\pi}(a) > \hat{\pi}(b)\}$$

255 *Kendall Tau.* Kendall's τ coefficient [42] is the normalized difference between
 the number of concordant, \mathcal{C} , and discordant pairs, \mathcal{D} :

$$\tau(\pi, \hat{\pi}) = \frac{\mathcal{C} - \mathcal{D}}{\frac{1}{2}k(k-1)}$$

where $\frac{1}{2}k(k-1)$ is the number of possible pairwise combinations, $\binom{k}{2}$. The
 values of this coefficient range from $[-1, 1]$, where $\tau(\pi, \pi) = 1$ (i.e. for equal
 rankings) and $\tau(\pi, \pi^{-1}) = -1$, where π^{-1} denotes the inverse order of π (e.g.
 260 $\pi = (1, 2, 3, 4)$ and $\pi^{-1} = (4, 3, 2, 1)$). Kendall's τ can also be computed in the
 presence of ties, using tau-b [43].

An alternative measure is the Spearman's rank correlation coefficient [44].

Gamma coefficient. If we want to measure the correlation between two partial
 orders (subrankings), or between total and partial orders, we can use the Gamma
 265 coefficient [45]:

$$\gamma(\pi, \hat{\pi}) = \frac{\mathcal{C} - \mathcal{D}}{\mathcal{C} + \mathcal{D}}$$

which is equivalent to Kendall's τ coefficient for strict total orders, because
 $\mathcal{C} + \mathcal{D} = \frac{1}{2}k(k-1)$.

Weighted rank correlation measures. When it is important to give more rele-
 vance to higher ranks, a weighted rank correlation coefficient can be used. They

270 are typically adaptations of existing similarity measures, such as ρ_w [46], which
is based on Spearman’s coefficient.

These correlation measures are not only used for evaluation estimation, they
can be used in the learning [7] and pre-processing [37] methods. Since Kendall’s
 τ has been used for evaluation in many recent Label Ranking studies [2, 36], we
275 use it here as well.

The accuracy of a label ranker can be estimated by averaging the values of any
of the measures explained here, over the rankings predicted for a set of test
examples. Given a dataset, $D = \{ \langle x_i, \pi_i \rangle \}$, $i = 1, \dots, n$, the usual resampling
280 strategies, such as holdout or cross-validation, can be used to estimate the
accuracy of a Label Ranking algorithm.

4. Label Ranking Association Rules

Association rules were originally proposed for descriptive purposes. However,
they have been adapted for predictive tasks such as classification (e.g., [6]).
285 Given that Label Ranking is a predictive task, the adaptation of AR for Label
Ranking comes in a natural way. A *Label Ranking Association Rule* (LRAR)
[7] is defined as:

$$A \rightarrow \pi$$

where $A \subseteq \text{desc}(\mathbb{X})$ and $\pi \in \Omega$. Let \mathcal{R}_π be the set of *Label Ranking association
rules* generated from a given dataset. When an instance x is covered by the
290 rule $A \rightarrow \pi$, the predicted ranking is π . A rule $r_\pi : A \rightarrow \pi, r_\pi \in \mathcal{R}_\pi$, covers an
instance x , if $A \subseteq \text{desc}(x)$.

We can use the CAR framework [6] for LRAR, by transforming each ranking
into a class. However this approach has two important problems. First, the
number of classes can be extremely large, up to a maximum of $k!$, where k is
295 the size of the set of labels, \mathcal{L} . This means that the amount of data required to
learn a reasonable mapping $\mathbb{X} \rightarrow \Omega$ is unreasonably large.

The second disadvantage is that this approach does not take into account the differences in nature between Label Rankings and classes. In classification, two examples either have the same class or not. In this regard, Label Ranking
 300 is more similar to regression than to classification. In regression, a large number of observations with a given target value, say 5.3, increases the probability of observing similar values, say 5.4 or 5.2, but not so much for very different values, say -3.1 or 100.2. This property must be taken into account in the induction of prediction models. A similar reasoning can be made in Label Ranking. Let us
 305 consider the case of a data set in which ranking $\pi_a = (1, 2, 3, 4)$ occurs in 1% of the examples. Treating rankings as classes would mean that $P(\pi_a) = 0.01$. Let us further consider that the rankings $\pi_b = (1, 2, 4, 3)$, $\pi_c = (1, 3, 2, 4)$ and $\pi_d = (2, 1, 3, 4)$, which are obtained from π_a by swapping a single pair of adjacent labels, occur in 50% of the examples. Taking into account the stochastic nature
 310 of these rankings [2], $P(\pi_a) = 0.01$ seems to underestimate the probability of observing π_a . In other words it is expected that the observation of π_b , π_c and π_d increases the probability of observing π_a and vice-versa, because they are similar to each other.

This affects even rankings which are not observed in the available data. For
 315 example, even though a ranking is not present in the dataset it would not be entirely unexpected to see it in future data. This also means that it is possible to compute the probability of unseen rankings.

To take all this into account, similarity-based interestingness measures were proposed to deal with rankings [7].

320 4.1. Interestingness measures in Label Ranking Association Rules

As mentioned before, because the degree of similarity between rankings can vary, similarity-based measures can be used to evaluate LRAR. These measures are able to distinguish rankings that are *very similar* from rankings that are *very distinct*. In practice, the measures described below can be applied to existing
 325 rule generation methods [7] (e.g. APRIORI [4]).

Support. The support of a ranking π should increase with the observation of similar rankings and that variation should be proportional to the similarity. Given a measure of similarity between rankings $s(\pi_a, \pi_b)$, we can adapt the concept of support of the rule $A \rightarrow \pi$ as follows:

$$\text{sup}_{lr}(A \rightarrow \pi) = \frac{\sum_{i:A \subseteq \text{desc}(x_i)} s(\pi_i, \pi)}{n}$$

330 Essentially, what we are doing is assigning a weight to each target ranking π_i in the training data where $A \subseteq \text{desc}(x_i)$. The weights represent the contribution of π_i to the probability that π may be observed. Some instances $x_i \in \mathbb{X}$ give a strong contribution to the support count (i.e., 1), while others will give a weaker or even no contribution at all.

335 Any function that measures the similarity between two rankings or permutations can be used, such as Kendall’s τ [42] or Spearman’s ρ [44]. The function used here is of the form:

$$s(\pi_a, \pi_b) = \begin{cases} s'(\pi_a, \pi_b) & \text{if } s'(\pi_a, \pi_b) \geq \theta \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where s' is a similarity function. This general form assumes that below a given threshold, θ , is not useful to discriminate between different rankings, as they are so different. This means that, the support sup_{lr} of $A \rightarrow \pi_a$ will be based 340 only on the items of the form $\langle A, \pi_b \rangle$, for all π_b where $s'(\pi_a, \pi_b) > \theta$.

Many functions can be used as s' . However, given that the loss function we aim to minimize is known beforehand, it makes sense to use it to measure the similarity between rankings. Therefore, we use Kendall’s τ as s' .

345 Concerning the threshold, given that anti-monotonicity can only be guaranteed with non-negative values [47], it implies that $\theta \geq 0$. Therefore we think that $\theta = 0$ is a reasonable default value, because it separates between the positive and negative correlation between rankings.

Table 1 shows an example of a Label Ranking dataset represented according 350 to this approach. Instance $\{\mathcal{A}_1 = L, \pi_3\}$ (TID=1) contributes to the support count of the rule $\mathcal{A}_1 = L \rightarrow \pi_3$ with 1, as expected. However, that same

Table 1: An example of a Label Ranking dataset. (TID = Transaction ID)

		π_1	π_2	π_3
TID	\mathcal{A}_1	(1, 3, 2)	(2, 1, 3)	(2, 3, 1)
1	L	0.33	0.00	1.00
2	L	0.00	1.00	0.00
3	L	1.00	0.00	0.33

instance, will also give a contribution of 0.33 to the support count of the rule $\mathcal{A}_1 = L \rightarrow \pi_1$, given the similarity between their rankings. On the other hand, no contribution to the support of the rule $\mathcal{A}_1 = L \rightarrow \pi_2$ is given, because these rankings are clearly different. This means that $sup_{lr}(\mathcal{A}_1 = L \rightarrow \pi_3) = \frac{1+0.33}{3}$.

Confidence. The confidence of a rule $A \rightarrow \pi$ comes in a natural way if we replace the classical measure of support with the similarity-based sup_{lr} .

$$conf_{lr}(A \rightarrow \pi) = \frac{sup_{lr}(A \rightarrow \pi)}{sup(A)}$$

Improvement. Improvement in association rule mining is defined as the smallest difference between the confidence of a rule and the confidence of all sub-rules sharing the same consequent (Section 2.3). In Label Ranking, it is not suitable to compare targets simply as equal or different, as explained earlier. Therefore, to implement pruning based on improvement for Label Ranking, some adaptation is required as well. Given that the relation between target values is different from the classification setting, as discussed earlier, we have to limit the comparison between rules with different consequents, if $s'(\pi, \pi') \geq \theta$.

Improvement for LRARs is defined as:

$$imp_{lr}(A \rightarrow \pi) = \min(conf_{lr}(A \rightarrow \pi) - conf_{lr}(A' \rightarrow \pi'))$$

for $\forall A' \subset A$, and $\forall (\pi, \pi')$ where $s'(\pi', \pi) \geq \theta$.

As an illustrative example, consider the two rules $r_1 : A_1 \rightarrow (1, 2, 3, 4)$ and $r_2 : A_2 \rightarrow (1, 2, 4, 3)$, where A_2 is a superset of A_1 , $A_1 \subset A_2$. If $s'((1, 2, 3, 4), (1, 2, 4, 3)) \geq \theta$ then r_2 will only be kept if, and only if, $conf(r_1) - conf(r_2) \geq minImp$.

Lift. This is a measure of the independence between the consequent and the antecedent of the rule [48]. The adaptation of *lift* for LRAR is straightforward since it only depends the concept of support, for which a version for LRAR already exists:

$$lift_{lr}(A \rightarrow \pi) = \frac{sup_{lr}(A \rightarrow \pi)}{sup(A) \cdot sup_{lr}(\pi)}$$

375 4.2. Generation of LRAR

Given the adaptations of the interestingness measures proposed, the task of learning LRAR can be defined essentially in the same way as the task of learning AR, i.e. to identify a set of LRAR which have a support and a confidence higher than the thresholds defined by the user. More formally, given a training set
 380 $D = \{\langle x_i, \pi_i \rangle\}, i = 1, \dots, n$, the algorithm aims to create a set of high accuracy rules $\mathcal{R}_\pi = \{r_\pi : A \rightarrow \pi\}$ to cover a test set $T = \{\langle x_j \rangle\}, j = 1, \dots, s$. If \mathcal{R}_π does not cover some $x_j \in T$, a *default ranking* (Section 4.3.1) is assigned to it.

4.2.1. Implementation of LRAR in CAREN

The association rule generator² we use is CAREN [49]. CAREN implements
 385 an association rule algorithm to derive rule-based prediction models, like CAR and LRAR. For Label Ranking datasets, CAREN derives association rules where the consequent is a complete ranking.

CAREN is specialized in generating association rules for predictive models and employs a bitwise depth-first frequent pattern mining algorithm. Rule
 390 pruning is performed using a Fisher exact test [49]. Like CMAR [50], CAREN is a rule-based algorithm rather than itemset-based. This means that, frequent itemsets are derived at the same time as rules are generated, whereas itemset-based algorithms carry out the two tasks in two separated steps.

Rule-based approaches allow for different pruning methods. For example, let
 395 us consider the rule $A \rightarrow \lambda$, where λ is the most frequent class in the examples covering A . If $sup(A \rightarrow \lambda) < minsup$ then there is no need to search for a

²<http://www4.di.uminho.pt/~pja/class/caren.html> (accessed 10.02.17)

superset of A , A^* , since any rule of the form $A^* \rightarrow \lambda, A \subset A^*$ cannot have a support higher than *minsup*.

CAREN generates significant rules [23]. Statistical significance of a rule is
 400 evaluated using a Fisher Exact Test by comparing its support to the support
 of its direct generalizations. The direct generalizations of a rule $A \rightarrow C$ are
 $\emptyset \rightarrow C$ and $(A \setminus \{a\}) \rightarrow C$ where a is a single item.

The final set of rules obtained define the Label Ranking prediction model,
 which we can also refer to as the *label ranker*.

405 CAREN also employs prediction for strict rankings using *consensus ranking*
 (Section 4.3), best rule, among others.

4.3. Prediction

A very straightforward method to generate predictions using a label ranker
 is used. The set of rules \mathcal{R}_π can be represented as an ordered list of rules, by
 410 some user-defined measure of relevance:

$$\langle r_{\pi_1}, r_{\pi_2}, \dots, r_{\pi_s} \rangle$$

As mentioned before, a rule $r_\pi^* : A^* \rightarrow \pi^*$ covers (or matches) an instance
 $x_i \in T$, if $A^* \subseteq desc(x_i)$. If only one rule, r_π^* , matches x_i , the predicted ranking
 for x_i is π^* . However, in practice, it is quite common to have more than one
 rule covering the same instance x_i , $\mathcal{R}_\pi^*(x_j) \subseteq \mathcal{R}_\pi$. In $\mathcal{R}_\pi^*(x_j)$ there can be rules
 415 with conflicting ranking recommendations. There simple approaches to address
 those conflicts, such as selecting the best rule, calculating the majority ranking,
 etc.

However, it has been shown that a ranking obtained by ordering the average
 ranks of the labels across all rankings minimizes the Spearman footrule distance
 420 to all those rankings [51]. In other words, it maximizes the similarity according
 to Spearman's ρ [44], and, consequently [52] Kendall's τ . This can be referred
 to as the *average ranking* [11].

Given any set of rankings $\{\pi_i\}$ ($i = 1, \dots, s$) with k labels, we compute the

average ranks as:

$$\bar{\pi}(j) = \frac{\sum_{i=1}^s \pi_i(j)}{s}, j = 1, \dots, k \quad (2)$$

425 The *average ranking* $\bar{\pi}$ can be obtained if we rank the values of $\bar{\pi}(j), j = 1, \dots, k$. A weighted version of this method can be obtained by using the *confidence* or *support* of the rules in $\mathcal{R}_\pi^*(x_j)$ as weights.

4.3.1. Default rules

As in classification, in some cases, the label ranker might not find any rule
 430 that covers a given instance x_j , so $\mathcal{R}_\pi^*(x_j) = \emptyset$. To avoid this, we need to define a *default rule*, r_\emptyset , which can be used in such cases:

$$\{\emptyset\} \rightarrow \text{default ranking}$$

A *default class* is also often used in classification tasks [53], which is usually the majority class of the training set D . In a similar way, we could define the majority ranking as our *default ranking*. However, some Label Ranking
 435 datasets have as many rankings as instances, making the majority ranking not so representative.

As mentioned before, the *average ranking* (Equation 2) of a set of rankings, minimizes the distance to all rankings in that set [51]. Hence we can use the *average ranking* of the target rankings in the training data as the *default ranking*.

440 4.4. Parameter tuning

Due to the intrinsic nature of each different dataset, or even of the pre-processing methods used to prepare the data (e.g., the discretization method), the *minsup/minconf* needed to obtain a rule set \mathcal{R}_π , that covers all the examples, may vary significantly [54]. The trivial solution would be, for example, to
 445 set *minconf* = 0 which would generate many rules, hence increasing the coverage. However, this rule would probably lead to a lot of uninteresting rules as well, as the model would overfit the data. Then, our goal is to obtain a rule set \mathcal{R}_π which gives maximal coverage while keeping high confidence rules.

Let us define M as the coverage of the model i.e. the coverage of the set of
 450 rules \mathcal{R}_π . Algorithm 1 represents a simple, heuristic method to determine the
 $minconf$ that obtains the rule set such that a certain minimal coverage, $minM$,
 is guaranteed.

Algorithm 1 Confidence tuning algorithm

Given $minsup$ and $step$
 $minconf = 100\%$
while $M < minM$ **do**
 $minconf = minconf - step$
 Run CAREN with $(minsup, minconf)$ and determine M
end while
return $minconf$

This procedure has the important advantage that it does not take into ac-
 count the accuracy of the rule sets generated, thus reducing the risk of overfit-
 455 ting.

5. Pairwise Association Rules

Association rules use a sets of descriptors to represent meaningful subsets
 of the data [55], hence providing an easy interpretation of the patterns mined.
 Due to the intuitive representation, since its first application for market bas-
 460 ket analysis [56], they have become very popular in data mining and machine
 learning tasks (Mining rankings [5], classification [6], Label Ranking [7], etc).

LRAR proved to be an effective predictive model [7], however they are de-
 signed to find complete rankings. Despite its similarity measures, which take
 into account ranking noise, they do not capture subranking patterns because
 465 they will always try to infer complete rankings. On the other hand, association
 rules were used to find patterns within rankings [5], but without relating them
 to the values of the independent variables.

In this work, we propose a decomposition method to look for meaning-
 ful associations between independent variables and preferences (in the form

470 of pairwise comparisons), the Pairwise Association Rules (PAR), which can be regarded as predictive or descriptive model. We define a PAR as:

$$A \rightarrow \{\lambda_a \succ \lambda_b \oplus \lambda_b \succ \lambda_a \oplus \lambda_a = \lambda_b \oplus \lambda_a \perp \lambda_b \mid \lambda_a, \lambda_b \in \mathcal{L}\}$$

where, as in the original AR paper [4], we allow rules with multiple items, not only in the antecedent but also in the consequent. In other words, PAR can also have multiple sets of pairwise comparisons in the consequent.

475 Similar to RPC (Section 3.1.3), we decompose the target rankings into pairwise comparisons. Therefore, PAR can be obtained from data with strict, partial and incomplete rankings³.

Contrary to LRAR, we use the same interestingness measures that are also used in typical AR approaches, instead of the similarity-based versions defined 480 for Label Ranking problems, i.e. *sup*, *conf*, etc. This allows PAR to filter out non-frequent/interesting patterns without the need to derive strict rankings. When methods cannot find interesting rules with enough pairwise comparisons to define a strict ranking, then it can abstain from making some choices and, thus, obtain partial rankings, subrankings or even with sets of disjoint pairwise 485 comparisons.

Abstention is used in machine learning to describe the option to not make a prediction when the confidence in the output of a model is insufficient. The simplest case is classification, where the model can abstain itself to make a decision [57]. In the Label Ranking task, a method that makes partial abstentions 490 was proposed in [30]. A similar reasoning is used here both for predictive and descriptive models. Partial abstentions also make sense in PAR. Hence, the decision to abstain on certain pairwise preferences is defined by interest measures, such as *minconf* or *lift*.

More formally, let us define $D = \{ \langle x_i, \pi_i \rangle \}, i = 1, \dots, n$ where π_i can be a 495 *complete ranking*, *partial ranking* or a *sub-ranking*. For each π of size k , we can

³To derive the PAR, we added a pairwise decomposition method to the CAREN [49] software.

extract up to h pairwise comparisons. We consider 4 possible outcomes for each pairwise comparison:

- $\lambda_a \succeq \lambda_b$
- $\lambda_b \succeq \lambda_a$
- 500 • $\lambda_a = \lambda_b$ (indifference)
- $\lambda_a \perp \lambda_b$ (incomparability)

As an example, a PAR can be of the form:

$$A \rightarrow \lambda_1 \succ \lambda_4 \wedge \lambda_3 \succ \lambda_1 \wedge \lambda_1 \perp \lambda_2$$

The consequent can be simplified into $\lambda_3 \succ \lambda_1 \succ \lambda_4$ or represented as a sub-ranking $\pi = (2, 0, 1, 3)$.

505 6. Experimental Results

In this section, we start by describing the datasets used in the experiments, then we introduce the experimental setup and finally present the results obtained.

6.1. Datasets

510 The Label Ranking datasets in this work (Table 2) were taken from the Data Repository of Paderborn University ⁴.

To illustrate domain-specific interpretations of the results, we experiment with two additional datasets. We use Algae [58], an adapted dataset from the 1999 COIL Competition [59], concerning the frequencies of algae populations in different environments ⁵. The original dataset consisted of 340 examples, each
515 representing measurements of a sample of water from different European rivers

⁴<https://www-old.cs.uni-paderborn.de/fachgebiete/intelligente-systeme/software/label-ranking-datasets.html> (accessed 10.02.17)

⁵<https://data.mendeley.com/datasets/spwmg2z7cv/> (accessed 10.02.17)

on different periods. The measurements include concentrations of chemical substances like nitrogen (in the form of nitrates, nitrites and ammonia), oxygen and chlorine. Also the pH, season, river size and its flow velocity were registered. For each sample, the frequencies of 7 types of algae were also measured. In this work, we considered the algae concentrations as preference relations by ordering them from larger to smaller concentrations. Those with 0 frequency are placed in last position and equal frequencies are represented with ties. Missing values in the independent variables were set to 0.

Finally, the Sushi preference dataset [12], which is composed of demographic data about 5000 people and sushi preferences, is also used. Each person sorted a set of 10 different sushi types by preference. The 10 types of sushi, are a) shrimp, b) sea eel, c) tuna, d) squid, e) sea urchin, f) salmon roe, g) egg h) fatty tuna, i) tuna roll and j) cucumber roll. Since the attribute names were not transformed in this dataset, it is particularly useful for the interpretation of the patterns extracted.

Table 2 also presents a simple measure of the diversity of the target rankings, the *Unique Ranking Proportion*, U_π . U_π is the proportion of distinct target rankings for a given dataset. As a practical example, the *iris* dataset has 5 distinct rankings for 150 instances, which results in $U_\pi = \frac{5}{150} \approx 3\%$.

6.2. Experimental setup

Continuous variables were discretized with two distinct methods: (1) *EDiRa* [37] and (2) *equal width* bins. *EDiRa* is the state of the art supervised discretization method in Label Ranking, while *equal width* is a simple, general method that serves as baseline.

The evaluation measure used in all experiments is Kendall’s τ (Section 3.2). A ten-fold cross-validation was used to estimate the value for each experiment. The generation of LRAR and PAR was performed with CAREN [49] which uses a depth-first approach.

The confidence tuning method described earlier (Algorithm 1) was used to set parameters. We consider that 5% seems a reasonable step value because the

Table 2: Summary of the datasets.

Datasets	#examples	#labels	#attributes	U_π
bodyfat	252	7	7	94%
calhousing	20,640	4	4	0.1%
cpu-small	8,192	5	6	1%
elevators	16,599	9	9	1%
fried	40,769	5	9	0.3%
glass	214	6	9	14%
housing	506	6	6	22%
iris	150	3	4	3%
segment	2310	7	18	6%
stock	950	5	5	5%
vehicle	846	4	18	2%
vowel	528	11	10	56%
wine	178	3	13	3%
wisconsin	194	16	16	100%
Algae	316	7	10	72%
Sushi	5000	10	10	98%

$minconf$ value can be found in, at most, 20 iterations. Given that a common value for the $minsup$ in association rule mining is 1%, we use it as default, except is stated otherwise. We define the $minM$ as 95%, to get a reasonable coverage, and $minImp = 1\%$, to avoid rule explosion.

In terms of similarity functions, we use a normalized Kendall τ between the interval $[0, 1]$ as our similarity function s' (Equation 1).

6.3. Results with LRAR

In the experiments described in this section, we analyze the performance of LRAR from different perspectives, namely *accuracy*, *number of rules* and *average confidence*, as the similarity threshold θ varies. We expect to understand

the impact of using similarity measures in the generation of LRAR and provide some insights about its usage.

LRAR, despite being based on similarity measures, are consistent with the
560 classical concepts underlying association rules. A special case is when $\theta = 1$,
where, as in CAR, only equal rankings are considered. Therefore, by varying the
threshold θ we also understand how similarity-based interest measures ($0 \leq \theta <$
1) contribute to the accuracy of the model, in comparison to frequency-based
approaches ($\theta = 1$).

565 We would also like to understand how some properties of the data relate the
sensitivity to θ . We can extract two simple measures of ranking diversity from
the datasets, the *Unique Ranking Proportion* (U_π), described earlier, and the
ranking entropy [37].

6.3.1. Sensitivity analysis: Accuracy

570 In Figure 1, we can see the behavior of the accuracy of CAREN varying
the value of θ . It shows that, in general, there is a tendency for the accuracy
to decrease as θ gets closer to 1. This happens in 12 out of the 14 datasets
analyzed. On the other hand, in 9 out of 14 datasets, the accuracy is rather
stable in the range $\theta \in [0, 0.6]$.

575 If we take into consideration that the model ignores the similarity between
rankings for $\theta = 1$, the results indicate that, as expected, there is advantage in
using the more flexible approach (i.e. taking ranking similarity into account)
compared to the strict classification approach (i.e. using CAR). Two extreme
cases are *fried* and *wisconsin*, where CAREN was not able to find any LRAR
580 for $\theta = 1$ ⁶.

Let us consider the *accuracy range*, the maximum accuracy minus the mini-
mum accuracy. To find out which datasets are more likely to be affected by the
choice of θ , we can compare their ranking entropy with the measured *accuracy*
range (In interest of space, we do not include the specific values here but they

⁶The *default rule* was not used in these experiments because it is not related to θ .

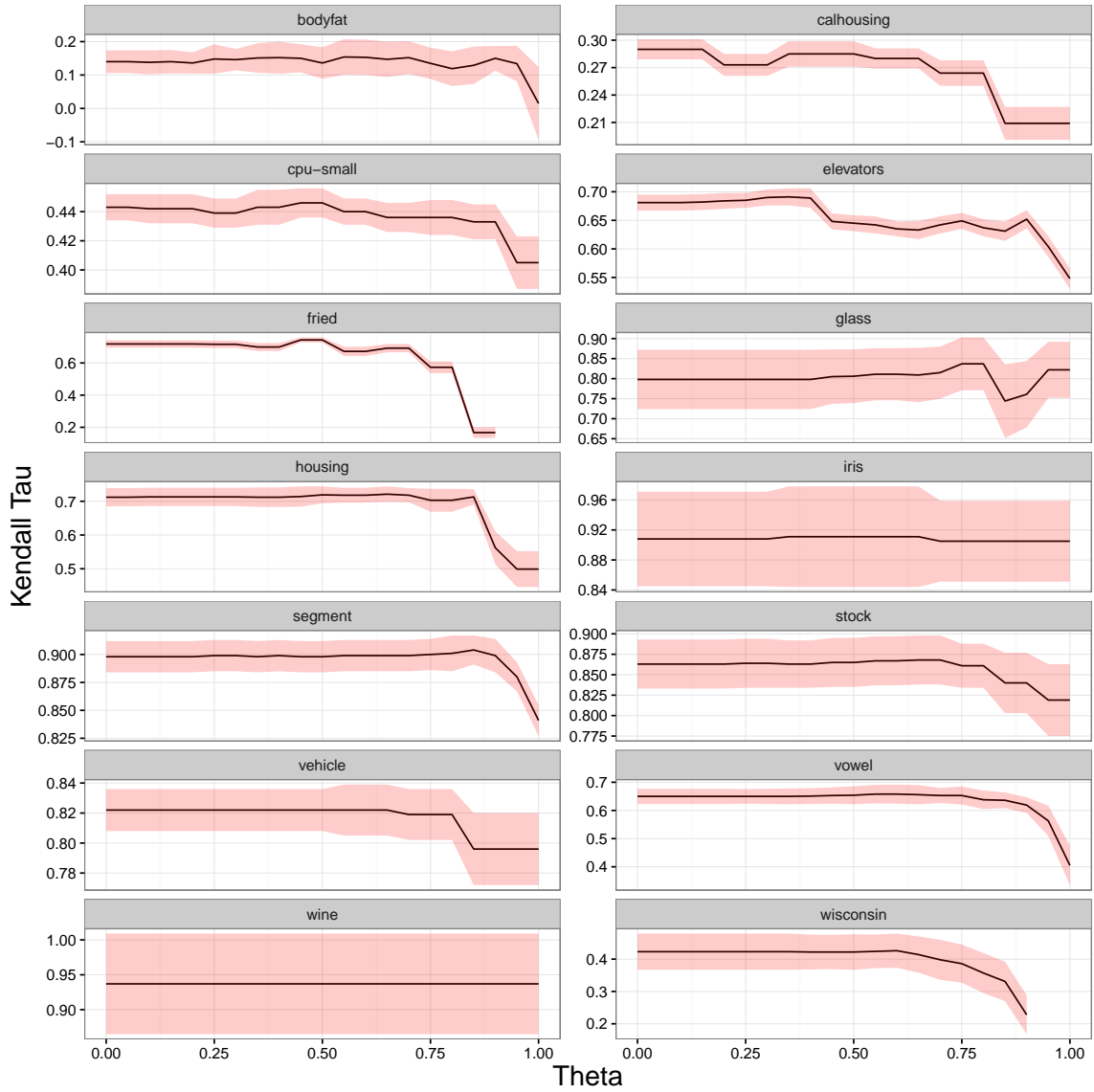


Figure 1: Average accuracy (Kendall τ) of CAREN as the value θ varies. (The shaded area represents the standard deviation)

585 can be easily estimated from Figure 1). In Figure 2, we compare the accuracy
 range with the *ranking entropy* [37]. We can see that, the higher the entropy,
 the more the accuracy can be affected by the choice of θ .

Results seem to indicate that, when mining LRAR in datasets with low
 ranking entropy, the choice of θ is not so relevant. On the other hand, as the
 590 entropy gets higher, reasonable values are in the range $0 \leq \theta \leq 0.6$.

Another interesting observation can be made regarding *fried*. Despite the
 fact that it has a very low proportion of unique rankings, $U_\pi(\textit{fried}) = 0.3\%$
 (Table 2) its entropy is quite high (Figure 2). For this reason, it makes it more
 sensitive to θ , as seen in Figure 1. On the other hand, *iris* and *wine*, with very
 595 low entropy, seem unaffected by θ .

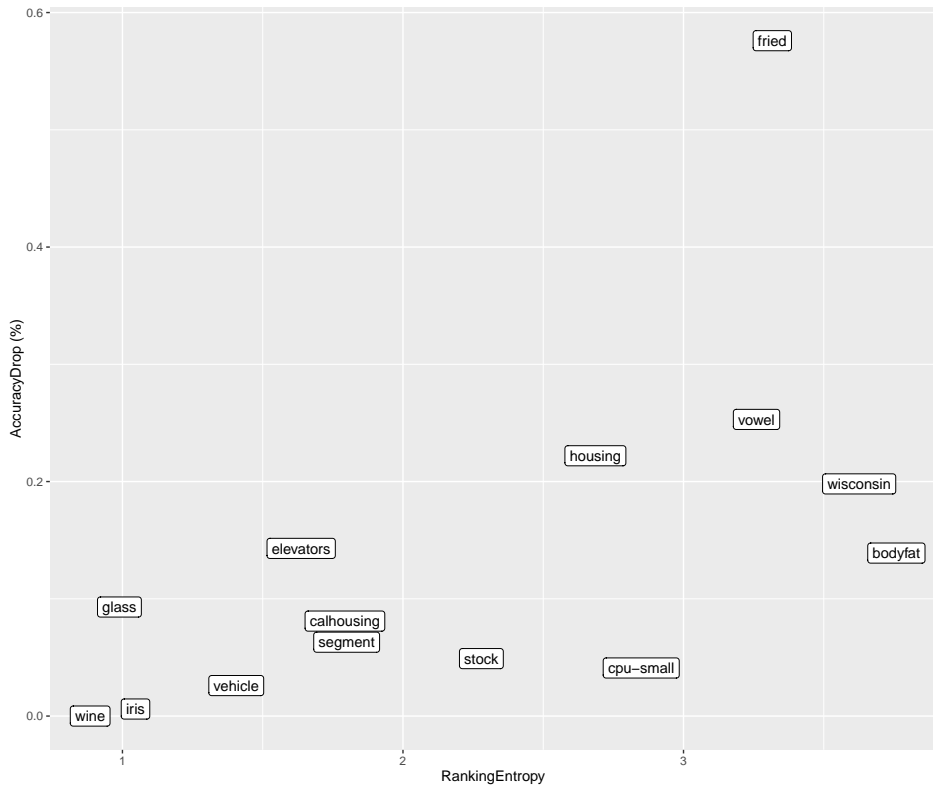


Figure 2: Accuracy range (Kendall τ) of CAREN in comparison to ranking entropy.

6.3.2. Sensitivity analysis: Number of rules

Ideally, we would like to obtain a small number of rules with high accuracy. However, such a balance is not expected to happen frequently. Ultimately, as accuracy is the most important evaluation criterion, if a reduction in the
600 number of rules comes with a high cost in accuracy, it is better to have more rules. Thus, it is important to understand how the number of LRAR varies with the similarity threshold θ , while taking the impact in the accuracy of the model into account as well.

In Figure 3, we see how many LRAR are generated per dataset as θ varies.
605 The majority of the plots, 10 out of 14, show a decrease in the number of rules as θ gets closer to 1. As discussed before, the accuracy in general also decreases as $\theta \geq 0.6$, so let us focus on $\theta \in [0, 0.6]$.

In the interval $\theta \in [0, 0.6]$, the number of rules generated is quite stable in 9 out of 14 datasets. In the first half of this interval, $\theta \in [0, 0.3]$, it is even more
610 remarkable for 13 datasets.

We expect the number of rules to decrease as θ increases, however, results show that the number of rules does not decrease so much, especially for values up to 0.3. This is due to the fact that θ is also used in the pruning step (Section 4.1), reducing the number of rules against which the improvement of an extension is
615 measured and, thus, increasing the probability of an extension not being kept in the model. This means that pruning is being effective in the reduction of LRAR. As mentioned before, $imp_{lr}(A \rightarrow \pi)$ not only compares rules $A' \rightarrow \pi$ where $A' \subset A$, but also rules $A \rightarrow \pi'$ where $S'(\pi', \pi) \geq \theta$. In other words, with the $minImp_{lr}$ we are pruning LRAR with similar rankings too.

620 These results do not lead to any strong conclusions about the ideal value for θ regarding the number of rules. However, they are in line with the previous analysis of *accuracy*.

6.3.3. Sensitivity analysis: Minimum Confidence

As described earlier, we use a greedy algorithm to automatically adjust the
625 minimum confidence in order to reduce the number of examples that are not

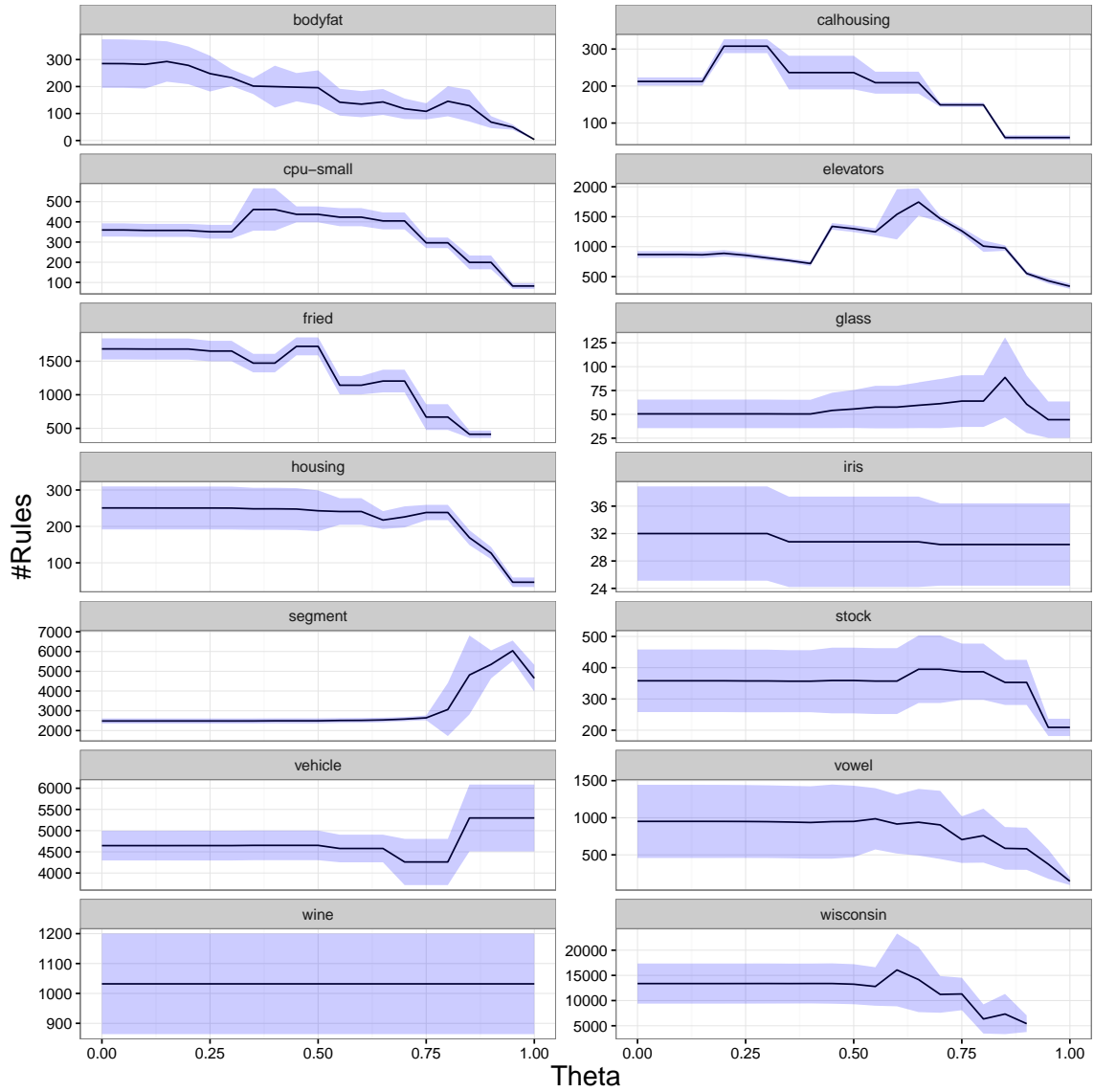


Figure 3: Number of Label Ranking Association Rules generated by CAREN as the value θ varies. (The shaded area represents the standard deviation)

covered by any rule. This means that different values of *minconf* depend on both the dataset and the value of θ , as seen in Figure 4.

In general, the *minconf* decreases in a monotonic way as θ increases. As $\theta \approx 1$ the *minconf* gets to its minimum on 13 out of 14 datasets, which is
630 consistent with the accuracy plots (Figure 1). This means that, if we want to generate rules with as much confidence, as measured by *minconf*, as possible, we should use the minimum θ , i.e. $\theta = 0$.

6.3.4. Sensitivity analysis: Support versus accuracy

We vary the minimum support threshold, *minsup*, to test how it affects the
635 accuracy of our learner. A similar study has been carried out on CBA [60]. Specifically, we vary the *minsup* from 0.1% to 10%, using a step size of 0.1%. Due to the complexity of these experiments, we only considered the six smallest datasets.

In general, as we increase *minsup* the accuracy decreases, which is a strong
640 indicator that the support should be small (Figure 5). All lines are monotonically decreasing, i.e. either the values remain constant or they decrease as *minsup* increases.

From a different perspective, the changes are generally very small for *minsup* \in [0.1%, 1.0%]. Considering that lower *minsup* generate potentially more rules, we
645 recommend *minsup* = 1% as a reasonable value to start experiments with.

Discretization techniques. To test the influence of the discretization method used, we compared *EDiRa* with a non-supervised discretization method, *equal width*.

In general, the accuracy had the same behavior, as a function of θ , as with
650 *EDiRa*, i.e. the results are highly correlated (Figure 6). However, the supervised approach is consistently better. These results add further evidence that *EDiRa* is a suitable discretization method for Label Ranking [37].

Similar behavior was observed concerning the number of rules generated and the minimum confidence, but are not presented here in interest of space.

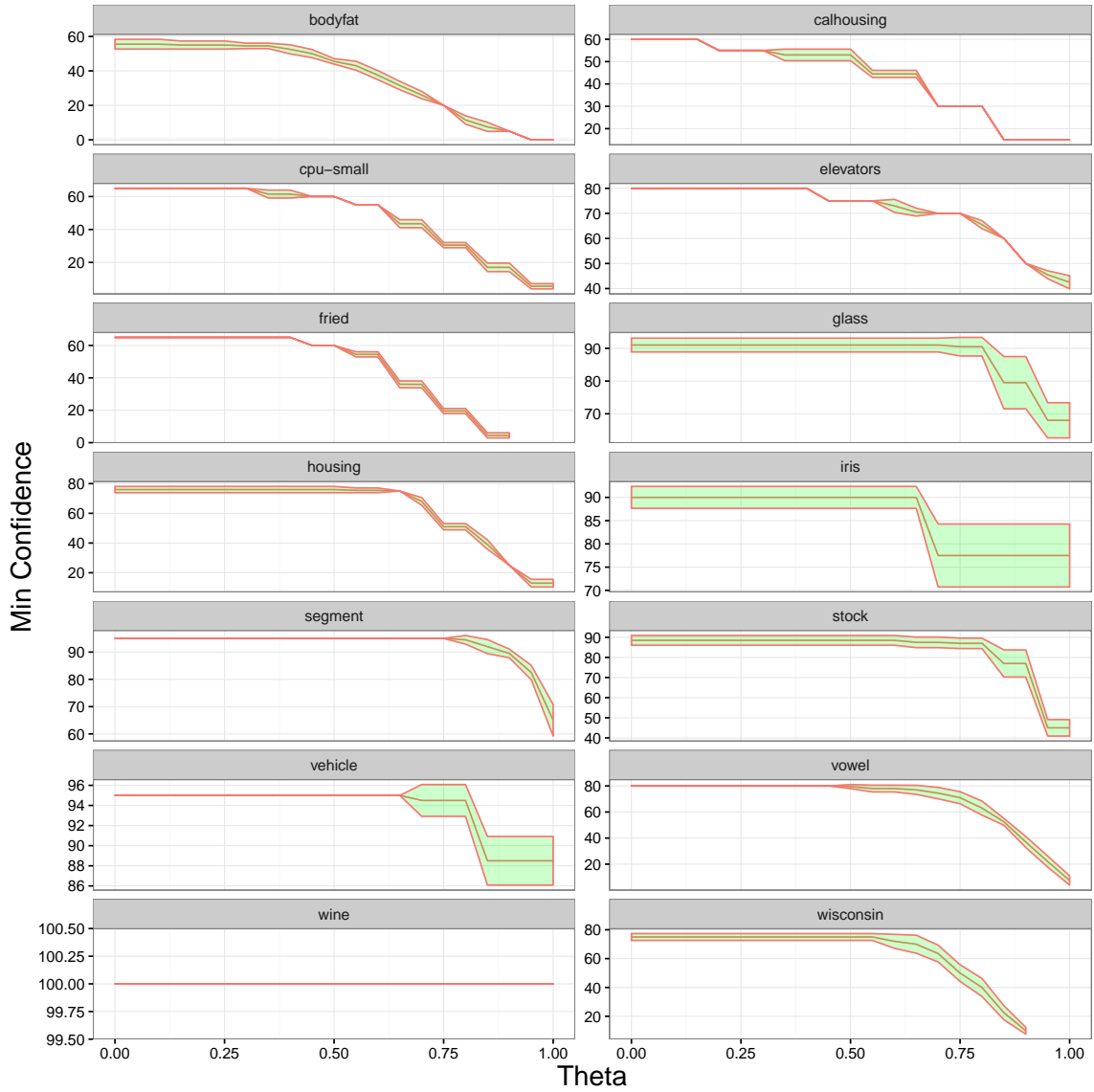


Figure 4: Minimum confidence as the value θ varies. (The shaded area represents the standard deviation)

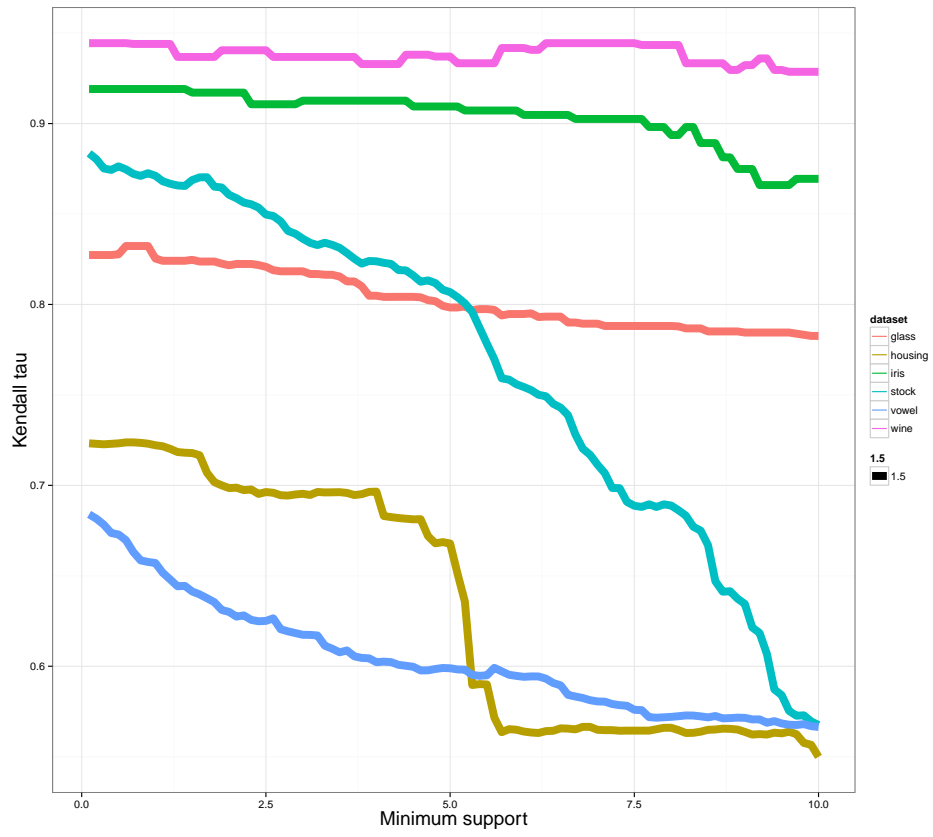


Figure 5: Average accuracy (Kendall τ) of CAREN as the value *minsup* varies.

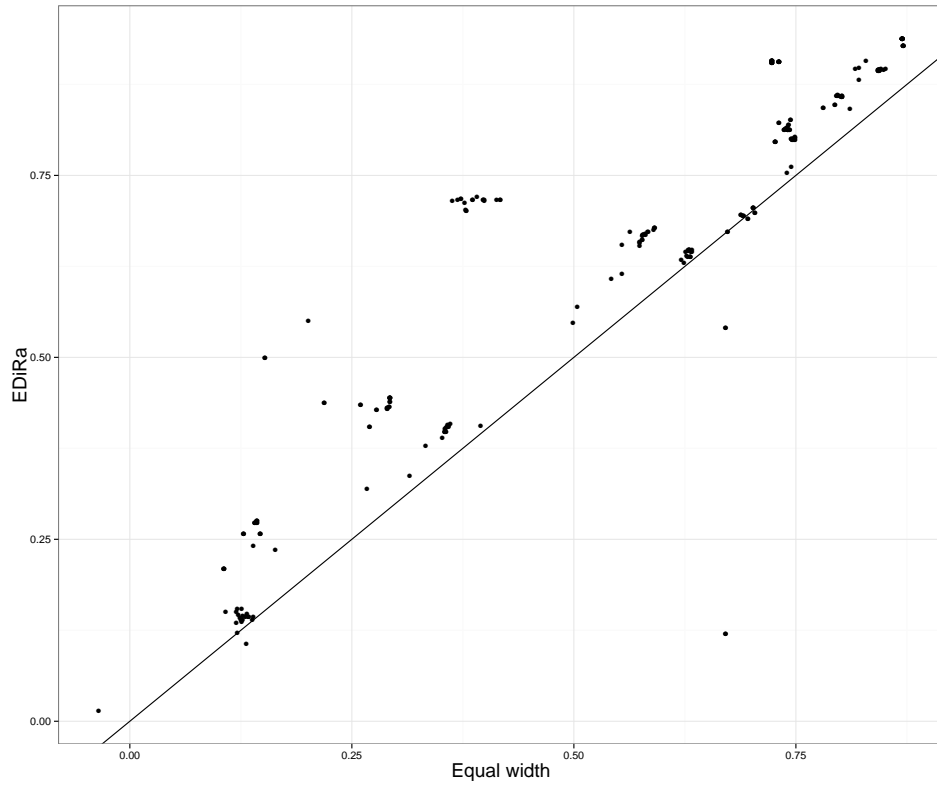


Figure 6: Ranking accuracy (Kendall τ) of CAREN after the discretization of data using *equal width* and *EDiRa*. This plot aggregates all the experiments carried out for each dataset, which means that each dataset is represented multiple times, with different parameter settings.

655 *Summary.* It is well known that general, simple rules to set parameters of machine learning algorithms do not exist. Nevertheless it is good to know where reasonable values lie. Hence, we think that $\theta \in [0.5, 0.6]$ and $minsup = 1\%$ are good default values for LRAR with CAREN. In terms of the discretization methods, our results confirm that a supervised approach, such as *EDiRa*, is a
660 good choice.

In Table 3 we compare the performance of LRAR with three state of the art approaches, RBLR, which is an alternative rule-based approach [8], IB-PL, an instance-based approach for Label Ranking [31] and Ranking by Pairwise Comparison [24]. We used the parameter values recommended earlier: the data
665 was discretized with the *EDiRa* method, θ was set to 0.5 and $minsup$ to 1%. It is important to note that the results presented for the other methods are the published results, we did not implement the mentioned approaches.

From Table 3, we see that LRAR are clearly a competitive approach, since their accuracy is in line with the reported values of other approaches. We can
670 conclude that LRAR are able to learn relevant patterns from Label Ranking data.

The lack of results for the RBLR and RPC on some datasets might be due to the size of the rankings in the training data. Both have a decomposition process that transforms the number of training examples into $n(k(k-1)/2)$ examples,
675 where n is the number of examples in the original data set and k the number of labels. Because of that the training time can increase dramatically [8].

6.4. Results with PAR

In this work, we use PAR as a descriptive model, to find patterns concerning subsets of labels. We focus in a descriptive task for two reasons. One is to make
680 the approach simpler and the other is to complement the predictive LRAR approach.

The minimum support and confidence presented here define the abstention level of the model. $minsup$ and $minconf$ were adjusted manually to generate a small enough set of rules to allow manual inspection (between 150 and 200).

Table 3: Results obtained on Label Ranking datasets using 4 different approaches (the mean accuracy is represented in terms of Kendall’s tau, τ).

	LRAR	RBLR	IB-PL	RPC
bodyfat	.136	-	.230	-
calhousing	.285	-	.326	-
cpu-small	.446	-	.495	-
elevators	.645	-	.721	-
fried	.743	-	.894	-
glass	.806	.882	.841	.882
housing	.719	-	.711	-
iris	.911	.956	.960	.885
segment	.898	-	.950	-
stock	.865	-	.922	-
vehicle	.822	.812	.859	.854
vowel	.654	.776	.851	.647
wine	.937	.883	.947	.921
wisconsin	.422	-	.479	-

685 Additionally, we set the minimum *lift* to 1.5. Despite that many interesting rules were found, due to space limitations we only present the most relevant.

Algae data. Using the Algae dataset, we found 179 PARs with *minsup* = 2 and *minconf* = 90. With *sup* = 2.2% and *conf* = 100%, the rule with the highest *lift* (approx. 6) was:

$$\begin{aligned} & \mathbf{Riversize} = \mathit{small} \wedge \mathbf{pH} \geq 37.9 \wedge \mathbf{Flowvelocity} = \mathit{high} \wedge \\ & \mathbf{Chloride} \geq 3.4 \wedge \mathbf{Nitrates\&Ammonia} \geq 18.5 \\ & \rightarrow L6 \succ L2 \wedge L5 \succ L7 \wedge L2 \succ L7 \end{aligned}$$

The consequent of this rule can be represented as $L6 \succ L2 \succ L7 \wedge L5 \succ L7$. Considering that the labels represent algae populations, this rule states that it is always true that, under these conditions, type 6 is more prevalent than type

690 2. It also states that type 7 is less prevalent than types 2, 5 and 6.

The rule with the second highest *lift* obtained, with *sup* = 3.1% and *conf* = 91%, is:

$$\begin{aligned} & \mathbf{Flowvelocity} = \mathit{medium} \wedge \mathbf{Nitrates\&Ammonia} < 18.5 \wedge \\ & \mathbf{Nitrogenasnitrates} < 7.9 \\ & \rightarrow L1 \succ L7 \wedge L7 \succ L3 \end{aligned}$$

The target of this rule is the partial ranking $L1 \succ L7 \succ L3$. If this PAR was used for prediction, the subranking $\pi = (1, 0, 3, 0, 0, 0, 2)$ would have been the prediction.

Sushi data. When analyzing the sushi dataset we got 166 rules with *minconf* = 70% and *minsup* = 1%. The following rule was found, with a *lift* of 1.95:

$$\begin{aligned} & \mathbf{Ageinterval} = 15 - 19 \wedge \mathbf{Sex} = \mathit{Male} \wedge \mathbf{Livedin} = \mathit{Eastern\ Japan} \\ & \rightarrow \mathbf{egg} \succ \mathbf{seaurchin} \wedge \mathbf{shrimp} \succ \mathbf{seaurchin} \end{aligned}$$

In the whole dataset, 37% of the people show this relative preferences $\mathbf{egg} \succ$
695 $\mathbf{seaurchin} \wedge \mathbf{shrimp} \succ \mathbf{seaurchin}$. This PAR shows that this number almost doubles (72%), if we consider males from Eastern Japan, aged between 15 – 19.

A related rule was also found concerning a different set of people, from a different age group and region ($sup = 1.1\%$, $conf = 71.6\%$ and $lift = 1.65$):

```

Ageinterval = 30 - 39  $\wedge$  Sex = Male  $\wedge$ 
Livesin = Western Japan  $\wedge$  Changedcity = Yes
 $\rightarrow$  seurchin  $\succ$  egg  $\wedge$ 
fattytuna  $\succ$  tunaroll  $\wedge$ 
tunaroll  $\succ$  cucumberroll  $\wedge$ 
fattytuna  $\succ$  egg

```

This rule includes one relative preference found in this group, **seurchin** \succ **egg**, which is the opposite to what was observed in the previous rule. Based on this information, we analyzed the data and found out that 75% of people that live
700 in Eastern Japan prefer **egg** to **seurchin** while 84% of people from Western Japan prefer **seurchin** to **egg**.

7. Conclusions

In this paper, we address the problem of finding association patterns in Label
Rankings. We present an extensive empirical analysis on the behavior of a Label
705 Ranking method, the CAREN implementation of Label Ranking Association
Rules. The performance was analyzed from different perspectives, *accuracy*,
number of rules and *average confidence*. The results show that, similarity-
based interest measures contribute positively to the accuracy of the model, in
comparison to frequency-based approaches, i.e. when $\theta = 1$.

710 The results confirm that LRAR are a viable Label Ranking tool which helps
solving complex Label Ranking problems (i.e. problems with high ranking en-
tropy). In comparison to other approaches, such as RPC, RBLR and IB-PL,
LRAR have the advantage to deliver interpretable results (in the form of asso-
ciation rules) and at the same time, without the need to decompose rankings,
715 which saves computational time. The results also enabled the identification of

some values for the parameters of the algorithm that can be used as default values.

Results also seem to indicate that, the higher the entropy, the more the accuracy can be affected by the choice of θ . The ranking entropy of a dataset
720 can be measured beforehand and the value of θ adjusted accordingly.

Additionally, we propose Pairwise Association Rules (PAR), which are association rules where the consequent represents multiple pairwise preferences. With PAR it is possible to obtain rules with complete, partial and incomplete rankings on the consequent. We illustrated the usefulness of this approach to
725 identify interesting patterns in Label Ranking datasets, which cannot be obtained with LRAR.

As future work, we will use PAR for predictive tasks.

Acknowledgments

This work is financed by the ERDF — European Regional Development
730 Fund through the Operational Programme for Competitiveness and Internationalization — COMPETE 2020 Programme within project POCI-01-0145-FEDER-006961, and by National Funds through the FCT — Fundação para a Ciência e a Tecnologia (Portuguese Foundation for Science and Technology) as part of project UID/EEA/50014/2013.

- 735 [1] J. Fürnkranz, E. Hüllermeier, Pairwise preference learning and ranking, in: Machine Learning: ECML 2003, 14th European Conference on Machine Learning, Cavtat-Dubrovnik, Croatia, September 22-26, 2003, Proceedings, 2003, pp. 145–156. doi:10.1007/978-3-540-39857-8_15.
- [2] W. Cheng, J. C. Huhn, E. Hüllermeier, Decision tree and instance-based
740 learning for label ranking, in: Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009, 2009, pp. 161–168. doi:10.1145/1553374.1553395.

- 745 [3] S. Vembu, T. Gärtner, Label ranking algorithms: A survey, in: Preference Learning., Springer, 2010, pp. 45–64. doi:10.1007/978-3-642-14125-6_3.
- [4] R. Agrawal, R. Srikant, Fast algorithms for mining association rules in large databases, in: VLDB'94, Proceedings of 20th International Conference on Very Large Data Bases, September 12-15, 1994, Santiago de Chile, Chile, 1994, pp. 487–499.
- 750 [5] S. Henzgen, E. Hüllermeier, Mining rank data, in: Discovery Science - 17th International Conference, DS 2014, Bled, Slovenia, October 8-10, 2014. Proceedings, 2014, pp. 123–134. doi:10.1007/978-3-319-11812-3_11.
- [6] B. Liu, W. Hsu, Y. Ma, Integrating classification and association rule mining, Knowledge Discovery and Data Mining (1998) 80–86.
- 755 [7] C. R. de Sá, C. Soares, A. M. Jorge, P. J. Azevedo, J. P. da Costa, Mining association rules for label ranking, in: Advances in Knowledge Discovery and Data Mining - 15th Pacific-Asia Conference, PAKDD 2011, Shenzhen, China, May 24-27, 2011, Proceedings, Part II, 2011, pp. 432–443. doi:10.1007/978-3-642-20847-8_36.
- 760 [8] M. Gurrieri, X. Siebert, P. Fortemps, S. Greco, R. Slowinski, Label ranking: A new rule-based label ranking method, in: Advances on Computational Intelligence - 14th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems, IPMU 2012, Catania, Italy, July 9-13, 2012. Proceedings, Part I, 2012, pp. 613–623. doi:10.1007/978-3-642-31709-5_62.
- 765 [9] S. Greco, B. Matarazzo, R. Slowinski, J. Stefanowski, An algorithm for induction of decision rules consistent with the dominance principle, in: Rough Sets and Current Trends in Computing, Second International Conference, RSCTC 2000 Banff, Canada, October 16-19, 2000, Revised Papers, 2000, pp. 304–313. doi:10.1007/3-540-45554-X_37.
- 770

- [10] L. Todorovski, H. Blockeel, S. Dzeroski, Ranking with predictive clustering trees, in: Machine Learning: ECML 2002, 13th European Conference on Machine Learning, Helsinki, Finland, August 19-23, 2002, Proceedings, 2002, pp. 444–455. doi:10.1007/3-540-36755-1_37.
- 775 [11] P. Brazdil, C. Soares, J. P. da Costa, Ranking learning algorithms: Using IBL and meta-learning on accuracy and time results, Machine Learning 50 (3) (2003) 251–277. doi:10.1023/A:1021713901879.
- [12] T. Kamishima, Nantonac collaborative filtering: recommendation based on order responses, in: Proceedings of the Ninth ACM SIGKDD International
780 Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 24 - 27, 2003, 2003, pp. 583–588. doi:10.1145/956750.956823.
- [13] R. Janicki, W. W. Koczkodaj, A weak order approach to group ranking, Comput. Math. Appl. 32 (2) (1996) 51–59. doi:10.1016/0898-1221(96)
785 00102-2.
URL [http://dx.doi.org/10.1016/0898-1221\(96\)00102-2](http://dx.doi.org/10.1016/0898-1221(96)00102-2)
- [14] Y. Zhang, D.-Y. Yeung, Multilabel relationship learning, ACM Trans. Knowl. Discov. Data 7 (2) (2013) 7:1–7:30. doi:10.1145/2499907.2499910.
790 URL <http://doi.acm.org/10.1145/2499907.2499910>
- [15] E. Omiecinski, Alternative interest measures for mining associations in databases, IEEE Trans. Knowl. Data Eng. 15 (1) (2003) 57–69. doi:10.1109/TKDE.2003.1161582.
- [16] M. Halkidi, M. Vazirgiannis, Quality assessment approaches in data mining,
795 in: Data Mining and Knowledge Discovery Handbook, 2nd ed., Springer, 2010, pp. 613–639. doi:10.1007/978-0-387-09823-4_31.
- [17] J. S. Park, M. Chen, P. S. Yu, An effective hash based algorithm for mining association rules, in: Proceedings of the 1995 ACM SIGMOD International

- 800 Conference on Management of Data, San Jose, California, May 22-25, 1995.,
1995, pp. 175–186. doi:10.1145/223784.223813.
- [18] S. Brin, R. Motwani, J. D. Ullman, S. Tsur, Dynamic itemset counting and
implication rules for market basket data, in: SIGMOD 1997, Proceedings
ACM SIGMOD International Conference on Management of Data, May
13-15, 1997, Tucson, Arizona, USA., 1997, pp. 255–264. doi:10.1145/
805 253260.253325.
- [19] J. S. Park, M. Chen, P. S. Yu, Efficient parallel and data mining for
association rules, in: CIKM '95, Proceedings of the 1995 International
Conference on Information and Knowledge Management, November 28 -
December 2, 1995, Baltimore, Maryland, USA, 1995, pp. 31–36. doi:
810 10.1145/221270.221320.
- [20] S. Thomas, S. Sarawagi, Mining generalized association rules and sequential
patterns using SQL queries, in: Proceedings of the Fourth International
Conference on Knowledge Discovery and Data Mining (KDD-98), New York
City, New York, USA, August 27-31, 1998, 1998, pp. 344–348.
- 815 [21] J. Han, J. Pei, Y. Yin, R. Mao, Mining frequent patterns without candidate
generation: A frequent-pattern tree approach, *Data Min. Knowl. Discov.*
8 (1) (2004) 53–87. doi:10.1023/B:DAMI.0000005258.31418.83.
- [22] R. J. B. Jr., R. Agrawal, D. Gunopulos, Constraint-based rule mining in
large, dense databases, *Data Min. Knowl. Discov.* 4 (2/3) (2000) 217–240.
820 doi:10.1023/A:1009895914772.
- [23] G. I. Webb, Discovering significant rules, in: Proceedings of the Twelfth
ACM SIGKDD International Conference on Knowledge Discovery and Data
Mining, Philadelphia, PA, USA, August 20-23, 2006, 2006, pp. 434–443.
doi:10.1145/1150402.1150451.
- 825 [24] E. Hüllermeier, J. Fürnkranz, W. Cheng, K. Brinker, Label ranking by

learning pairwise preferences, *Artif. Intell.* 172 (16-17) (2008) 1897–1916.
doi:10.1016/j.artint.2008.08.002.

[25] V. Chankong, Y. Haimes, *Multiobjective Decision Making: Theory and Methodology*, Dover Books on Engineering, Dover Publications, 2008.

830 URL <https://books.google.pt/books?id=o371DAAAQBAJ>

[26] J. Chomicki, Preference formulas in relational queries, *ACM Trans. Database Syst.* 28 (4) (2003) 427–466. doi:10.1145/958942.958946.

URL <http://doi.acm.org/10.1145/958942.958946>

[27] J. Fürnkranz, E. Hüllermeier, Preference learning: An introduction, in: *Preference Learning.*, Springer, 2010, pp. 1–17. doi:10.1007/978-3-642-14125-6_1.

835

[28] F. Brandenburg, A. Gleißner, A. Hofmeier, Comparing and aggregating partial orders with kendall tau distances, *Discrete Math., Alg. and Appl.* 5 (2). doi:10.1142/S1793830913600033.

[29] K. Brinker, E. Hüllermeier, Label ranking in case-based reasoning, in: *Case-Based Reasoning Research and Development, 7th International Conference on Case-Based Reasoning, ICCBR 2007, Belfast, Northern Ireland, UK, August 13-16, 2007, Proceedings, 2007*, pp. 77–91. doi:10.1007/978-3-540-74141-1_6.

840

URL http://dx.doi.org/10.1007/978-3-540-74141-1_6

845

[30] W. Cheng, M. Rademaker, B. D. Baets, E. Hüllermeier, Predicting partial orders: Ranking with abstention, in: *Machine Learning and Knowledge Discovery in Databases, European Conference, ECML PKDD 2010, Barcelona, Spain, September 20-24, 2010, Proceedings, Part I, 2010*, pp. 215–230. doi:10.1007/978-3-642-15880-3_20.

850

[31] W. Cheng, K. Dembczynski, E. Hüllermeier, Label ranking methods based on the plackett-luce model, in: *Proceedings of the 27th International Con-*

ference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel, 2010, pp. 215–222.

- 855 [32] S. Har-Peled, D. Roth, D. Zimak, Constraint classification: a new approach to multiclass classification, in: Proc. of the International Workshop on Algorithmic Learning Theory (ALT), Springer-Verlag, 2002, pp. 135–150.
- [33] S. Thrun, L. K. Saul, B. Schölkopf (Eds.), Advances in Neural Information Processing Systems 16 [Neural Information Processing Systems, NIPS 2003, 860 December 8-13, 2003, Vancouver and Whistler, British Columbia, Canada], MIT Press, 2004.
- [34] G. Lebanon, J. D. Lafferty, Conditional models on the ranking poset, in: Advances in Neural Information Processing Systems 15 [Neural Information Processing Systems, NIPS 2002, December 9-14, 2002, Vancouver, British 865 Columbia, Canada], 2002, pp. 415–422.
- [35] A. Aiguzhinov, C. Soares, A. P. Serra, A similarity-based adaptation of naive bayes for label ranking: Application to the metalearning problem of algorithm recommendation, in: Discovery Science - 13th International Conference, DS 2010, Canberra, Australia, October 6-8, 2010. Proceedings, 870 2010, pp. 16–26. doi:10.1007/978-3-642-16184-1_2.
- [36] C. R. de Sá, C. Soares, A. Knobbe, P. J. Azevedo, A. M. Jorge, Multi-interval discretization of continuous attributes for label ranking, in: Discovery Science - 16th International Conference, DS 2013, Singapore, October 6-9, 2013. Proceedings, 2013, pp. 155–169. doi:10.1007/ 875 978-3-642-40897-7_11.
- [37] C. R. de Sá, C. Soares, A. Knobbe, Entropy-based discretization methods for ranking data, Inf. Sci. 329 (2016) 921–936. doi:10.1016/j.ins.2015.04.022.
- [38] J. Dougherty, R. Kohavi, M. Sahami, Supervised and unsupervised discretization of continuous features, in: Machine Learning, Proceedings of 880

the Twelfth International Conference on Machine Learning, Tahoe City, California, USA, July 9-12, 1995, 1995, pp. 194–202.

- [39] W. Cheng, Label ranking with probabilistic models, Ph.D. thesis, University of Marburg (2012).
885 URL <http://archiv.ub.uni-marburg.de/diss/z2012/0493>
- [40] J. C. Fodor, M. R. Roubens, Fuzzy preference modelling and multicriteria decision support, Vol. 14 of Theory and Decision Library D., Springer Netherlands, Dordrecht, 1994. doi:10.1007/978-94-017-1648-2.
- [41] J. Fürnkranz, E. Hüllermeier (Eds.), Preference Learning, Springer, 2010.
890 doi:10.1007/978-3-642-14125-6.
- [42] M. Kendall, J. Gibbons, Rank correlation methods, Griffin London, 1970.
- [43] A. Agresti, Analysis of ordinal categorical data, Wiley series in probability and mathematical statistics, J. Wiley, Hoboken, 2010.
- [44] C. Spearman, The proof and measurement of association between two
895 things, *American Journal of Psychology* 15 (1904) 72–101.
- [45] W. H. K. Leo A. Goodman, Measures of association for cross classifications, *Journal of the American Statistical Association* 49 (268) (1954) 732–764.
- [46] J. Pinto da Costa, C. Soares, A weighted rank measure of correlation, *Australian and New Zealand Journal of Statistics* 47 (4) (2005) 515–529.
900 doi:10.1111/j.1467-842X.2005.00413.x.
- [47] J. Pei, J. Han, L. V. S. Lakshmanan, Mining frequent item sets with convertible constraints, in: *Proceedings of the 17th International Conference on Data Engineering*, April 2-6, 2001, Heidelberg, Germany, 2001, pp. 433–442. doi:10.1109/ICDE.2001.914856.
- 905 [48] P. J. Azevedo, A. M. Jorge, Comparing rule measures for predictive association rules, in: *Machine Learning: ECML 2007, 18th European Conference*

on Machine Learning, Warsaw, Poland, September 17-21, 2007, Proceedings, 2007, pp. 510–517. doi:10.1007/978-3-540-74958-5_47.

- [49] P. J. Azevedo, A. M. Jorge, Ensembles of jittered association rule classifiers, *Data Min. Knowl. Discov.* 21 (1) (2010) 91–129. doi:10.1007/s10618-010-0173-y. 910
- [50] W. Li, J. Han, J. Pei, CMAR: accurate and efficient classification based on multiple class-association rules, in: Proceedings of the 2001 IEEE International Conference on Data Mining, 29 November - 2 December 2001, San Jose, California, USA, 2001, pp. 369–376. doi:10.1109/ICDM.2001.989541. 915
- [51] J. Kemeny, J. Snell, *Mathematical Models in the Social Sciences*, MIT Press, 1972.
- [52] L. Winner, Nascar winston cup race results for 1975-2003, *Journal of Statistics Education: An international journal on the teaching and learning of statistics* 14 (3). 920
- [53] J. Han, M. Kamber, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2000.
- [54] B. Liu, W. Hsu, Y. Ma, Mining association rules with multiple minimum supports, in: Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Diego, CA, USA, August 15-18, 1999, 1999, pp. 337–341. doi:10.1145/312129.312274. 925
- [55] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer New York, New York, NY, 2009, Ch. Unsupervised Learning, pp. 485–585. doi:10.1007/978-0-387-84858-7_14. 930
- [56] R. Agrawal, T. Imielinski, A. N. Swami, Mining association rules between sets of items in large databases, in: Proceedings of the 1993 ACM SIGMOD

- International Conference on Management of Data, Washington, D.C., May
935 26-28, 1993., 1993, pp. 207–216. doi:10.1145/170035.170072.
- [57] P. L. Bartlett, M. H. Wegkamp, Classification with a reject option using
a hinge loss, *Journal of Machine Learning Research* 9 (2008) 1823–1840.
doi:10.1145/1390681.1442792.
- [58] C. R. de Sá, W. Duivesteijn, C. Soares, A. Knobbe, Exceptional pref-
940 erences mining, in: *Discovery Science*, 2016, p. 1–16. doi:10.1007/
978-3-319-46307-0_1.
- [59] K. Bache, M. Lichman, *UCI machine learning repository* (2013).
- [60] M. Iqbal, I. Mukhlash, H. M. Astuti, The comparison of cba algorithm and
cbs algorithm for meteorological data classification, *ISICO* 2013.