

Nikolaj Tatti. 2008. Maximum entropy based significance of itemsets. Knowledge and Information Systems (KAIS), accepted for publication.

© 2008 Springer Science+Business Media

Preprinted with kind permission of Springer Science and Business Media.

# Maximum entropy based significance of itemsets

Nikolaj Tatti

Received: 7 December 2007 / Accepted: 29 January 2008  
© Springer-Verlag London Limited 2008

**Abstract** We consider the problem of defining the significance of an itemset. We say that the itemset is significant if we are surprised by its frequency when compared to the frequencies of its sub-itemsets. In other words, we estimate the frequency of the itemset from the frequencies of its sub-itemsets and compute the deviation between the real value and the estimate. For the estimation we use Maximum Entropy and for measuring the deviation we use Kullback–Leibler divergence. A major advantage compared to the previous methods is that we are able to use richer models whereas the previous approaches only measure the deviation from the independence model. We show that our measure of significance goes to zero for derivable itemsets and that we can use the rank as a statistical test. Our empirical results demonstrate that for our real datasets the independence assumption is too strong but applying more flexible models leads to good results.

**Keywords** Binary data mining · Itemsets · Maximum entropy

## 1 Introduction

How significant is a given itemset? Itemsets are popular and well-studied patterns in binary data mining. The major drawback is that, given a dataset, there are exponential number of itemsets. Hence, we need to rank itemsets in order to prune the uninteresting ones.

Traditionally, the frequency of an itemset is used as a rank measure. The higher the frequency, the more significant is the itemset. Frequency has many virtues: It is easy to interpret and because of its property of anti-monotonicity there exist efficient algorithms for finding all frequent itemsets [2,3]. There are, however, major drawbacks. First, a frequent

---

A preliminary version appeared as “Maximum Entropy Based Significance of Itemsets”, In Proceedings of Seventh IEEE International Conference on Data Mining (ICDM 2007), pp 312–321, 2006 [32].

---

N. Tatti (✉)  
HIIT Basic Research Unit, Department of Computer Science,  
Helsinki University of Technology, Helsinki, Finland  
e-mail: ntatti@cc.hut.fi

itemset may be insignificant: An itemset  $AB$  may be frequent just because itemsets  $A$  and  $B$  are frequent. Second, an infrequent itemset may be significant: If itemsets  $A$  and  $B$  are frequent, the infrequency of  $AB$  is interesting information.

Alternative methods for ranking itemsets are suggested in [1, 6, 14]. These methods are discussed in more detail in Sect. 4. A common feature to these methods is that they compare the frequency of an itemset to an estimate obtained from the independence model. That is, the more the itemset deviates from the independence model, the more surprising, and thus the more significant, the itemset is.

Our proposal for ranking itemsets resembles the aforementioned approaches. We estimate the frequency of a given itemset from the frequencies of some selected sub-itemsets. Namely, we use Maximum Entropy for the estimation. This approach is more flexible than the independence model, since the independence model uses only the margins (the frequencies of itemsets of size 1) for prediction whereas our approach allows to use the information available from the itemsets of larger size. While our ranking method is based on well-known tools, no similar framework has been suggested previously.

Unlike the frequency, our measure is not decreasing with respect to set inclusion. Hence we cannot mine significant itemsets in a level-wise fashion. However, it turns out that in some cases we can prune a large set of uninteresting itemsets (w.r.t. the measure). Namely, if the itemset is derivable [8], then the measure is equal to 0. We also point out that can be used as a statistical test, thus providing a clear interpretation for the measure.

The rest of the paper is organized as follows: Preliminaries are given in Sect. 2. The definition and the properties of the measure are given in Sect. 3. We present related work in Sect. 4. Section 5 is devoted to experiments and finally we provide conclusions in Sect. 6.

## 2 Preliminaries and notation

In this section we review briefly theory of itemsets and also introduce some notation that will be used later on.

A *binary dataset*  $D$  is a collection of  $M$  binary vectors, *transactions*, having length  $K$ . Such dataset can be naturally represented as a matrix of size  $M \times K$ . We denote the number of transactions by  $|D| = M$ . To each column of the matrix we assign an *attribute*  $a_i$ . Let  $A = \{a_1, \dots, a_K\}$  be the collection of all attributes. An itemset  $X \subseteq A$  is a set of attributes.

We say that a transaction (binary vector)  $\omega$  *covers* an itemset  $X$  if  $a_i \in X$  implies  $\omega_i = 1$ . Given a dataset  $D$ , a *frequency* of an itemset  $X$  is a proportion of the transaction in  $D$  covering  $X$ . Note that if an itemset  $Y$  is a subset of  $X$ , then the frequency of  $Y$  is larger than or equal to the frequency of  $X$ . In other words, frequency is decreasing with respect to set inclusion.

A sample space  $\Omega$  is the set of all binary vectors of length  $K$ . We take a simplistic approach in defining distributions: A distribution  $p : \Omega \rightarrow [0, 1]$  is a function from a sample space  $\Omega$  to a real number between 0 and 1 such that  $\sum_{\omega \in \Omega} p(\omega) = 1$ . Given an itemset  $X$ , a frequency of  $X$  calculated from a distribution  $p$  is the probability of binary vector covering  $X$ . We denote this by

$$p(X = 1) = p(\omega \text{ covers } X).$$

A family of itemsets  $\mathcal{F}$  is called *anti-monotonic* or *downward closed* if every subset of each member of  $\mathcal{F}$  is also a member of  $\mathcal{F}$ . Note that a collection of  $\sigma$ -frequent itemsets, that is, itemsets having frequency larger than some given threshold  $\sigma$ , is downward closed. We are interested in three particular families:

- $\mathcal{I}$ , the family containing only itemsets of size 1.
- $\mathcal{C}$ , the family containing itemsets of size 1 and 2.
- $\mathcal{A}$ , the family containing all itemsets.

A *negative border*  $\text{negbord}(\mathcal{F})$  of the downward closed family  $\mathcal{F}$  is the set of itemsets just above  $\mathcal{F}$ . In other words,  $X \notin \mathcal{F}$  is member of  $\text{negbord}(\mathcal{F})$  if there is no proper subset  $Y \subset X$  such that  $Y \notin \mathcal{F}$ .

Given a dataset  $D$ , we say that an itemset  $X$  is *derivable* if by knowing the frequencies (calculated from  $D$ ) of each proper subset of  $X$  we can deduce the frequency of  $X$ . For example, if some subset of  $X$  has a frequency 0, then we know that  $X$  must also have frequency 0. Thus, in this case,  $X$  is derivable. An itemset that is not derivable is called *non-derivable*. A family of all non-derivable itemsets is downward closed [8].

### 3 Maximum entropy ranking

In this section we introduce our ranking method and discuss its theoretical properties. The fundamental idea behind our approach is to measure how surprising an itemset is compared to its subsets. In other words, we estimate the itemset frequency by using the frequencies of its subsets and compare how close is our estimation to the actual value. The estimation is done using maximum entropy method and the comparison is done using Kullback–Leibler divergence.

#### 3.1 Definition

Let  $D$  be a binary dataset and let  $\{a_1, \dots, a_K\}$  be its attributes. The number of columns in  $D$  is  $K$ . Assume that we are given  $G$ , an itemset we wish to rank. We define a projected dataset  $D_G$  by keeping only the attributes included in  $G$ .

Let  $\Omega_G = \{0, 1\}^{|G|}$  be a space of binary vectors of length  $|G|$ . We define an *empirical distribution*  $q_G : \Omega_G \rightarrow [0, 1]$  to be

$$q_G(\omega) = \frac{\text{Number of samples in } D_G \text{ equal to } \omega}{|D_G|}.$$

Our goal is to compare the distribution  $q_G$  to a distribution obtained by using maximum entropy [23], a method that we will describe next.

Assume now that we are given a family of itemsets  $\mathcal{F} \subseteq \mathcal{A}$  and let  $\theta_X$  be the frequency of  $X \in \mathcal{F}$  calculated from  $D$ . Our next step is to define an approximative distribution using only the itemsets in  $\mathcal{F}$ . In defining  $q_G$  we projected out the attributes outside  $G$ . Similarly, we are only interested in subsets of  $G$ . Hence we define a *projected family*  $\mathcal{F}_G$  to be

$$\mathcal{F}_G = \{X \in \mathcal{F} \mid X \subset G, X \neq G, X \neq \emptyset\}.$$

Note that  $\mathcal{F}_G$  may contain  $2^{|G|} - 2$  itemsets, at maximum. This is the case if  $\mathcal{F} = \mathcal{A}$ .

We say that a distribution  $p : \Omega_G \rightarrow [0, 1]$  *satisfies the itemsets*  $\mathcal{F}_G$  if for each itemset  $X \in \mathcal{F}_G$  and its frequency  $\theta_X$  we have

$$p(X = 1) = \theta_X.$$

Let  $\mathbb{P}$  be the set of all distributions satisfying the itemsets  $\mathcal{F}_G$ . This set is not empty since  $q_G \in \mathbb{P}$ . We select the distribution from  $\mathbb{P}$  maximizing the entropy

$$H(p) = - \sum_{\omega \in \Omega_G} p(\omega) \log p(\omega).$$

We denote this distribution by  $p^*$ . Note that  $p^*$  depends on  $G, \mathcal{F}$ , and  $\theta$  but we have omitted these variables from the notation for the sake of clarity.

We define the rank measure  $r(G; \mathcal{F}, D)$  to be the divergence between  $q_G$  and  $p^*$ , that is,

$$r(G; \mathcal{F}, D) = \sum_{\omega \in \Omega_G} q_G(\omega) \log \frac{q_G(\omega)}{p^*(\omega)}.$$

We omit  $D$  from the notation when the dataset is clear from the context.

*Example 1* Assume the simplest case where  $G = a$  is an itemset of size 1. Let  $\theta_G$  be the frequency of  $G$ . Note that  $\mathcal{F}_G = \emptyset$ , hence there are no constraints on selecting  $p^*$ . This means that  $p^*$  is the uniform distribution, that is,  $p^*(0) = p^*(1) = 1/2$ . In this case the measure is

$$r(a; \mathcal{F}) = (1 - \theta_G) \log(2(1 - \theta_G)) + \theta_G \log(2\theta_G)$$

obtaining its minimum when  $\theta_G = 1/2$  and is at its maximum when  $\theta_G = 0$  or  $\theta_G = 1$ .

We are mainly interested in three kinds of measures: The first is  $r(G; \mathcal{I})$  in which  $\mathcal{I}$  is the family of itemsets of size 1. In this case the Maximum Entropy distribution is equal to the independence model.

The second case is  $r(G; \mathcal{C})$ , where  $\mathcal{C}$  contains the itemsets of size 1 and 2. We can show that there exists a matrix  $B$  (see [11]) such that for the non-zero entries of  $p^*$  we have

$$p^*(\omega) \propto \exp(\omega^T B \omega).$$

Hence,  $r(G; \mathcal{C})$  can be seen as the measure of the deviation from the discrete Gaussian model.

Our third type of measure is  $r(G; \mathcal{A})$  in which  $p^*$  is predicted from all the proper sub-itemsets of  $G$ . In this case we can prove that for a certain set of real numbers  $r_i$  we have for the non-zero entries of  $p^*$

$$p^*(\omega) \propto \prod_{X_i \in \mathcal{A}_G} \exp(r_i I(\omega \text{ covers } X_i)),$$

where  $I$  is the indicator function [11]. We discuss the evaluation of our approach in Sect. 3.4.

### 3.2 Properties

In this section we discuss various properties of  $r(G)$ . We will first point the connection between  $r(G)$  and derivable itemsets and then discuss the use of  $r(G)$  as a statistical test.

**Theorem 2** *Let  $G$  be a derivable itemset. Then*

$$r(G; \mathcal{A}) = 0.$$

*Proof* We can argue that if we know the frequencies of all sub-itemsets of  $G$ , we can derive the distribution  $q_G$  and vice versa. This implies that there is one-to-one correspondence between the distribution  $p \in \mathbb{P}$  satisfying the itemsets  $\mathcal{A}_G$  and the frequency  $p(G = 1)$ . Since we can derive the frequency of  $G$  from  $\mathcal{A}_G$ , it follows that  $\mathbb{P} = \{q_G\}$ , and hence  $p^* = q_G$ .  $\square$

We can reformulate the previous theorem in a stronger form by pointing out that we need to know only non-derivable itemsets.

**Theorem 3** *Let  $\mathcal{F}$  be a family of all non-derivable itemsets. Let  $G$  be outside of  $\mathcal{F}$ . Then  $r(G; \mathcal{F}) = 0$ .*

*Proof* Since all unknown sub-itemsets of  $G$  are derivable from  $\mathcal{F}_G$ , the argument of Theorem 2 holds.  $\square$

The following theorem provides the interpretation to the value of  $r(G)$  and points out that we can use  $r(G)$  as a statistical test.

**Theorem 4** *Let  $G$  be a non-derivable itemset. Under the 0-hypothesis that  $G$  is distributed according to  $p^*$ , the quantity  $2|D| r(G; \mathcal{A})$  is distributed asymptotically as  $\chi^2$  with degree 1 of freedom.*

Theorem 4 is a special case of the following more general statement.

**Theorem 5** *Let  $G$  be a non-derivable itemset and let  $\mathcal{F}$  be an itemset family. Define  $\mathcal{H}$  to be*

$$\mathcal{H} = \{X \in \mathcal{A} \mid X \subseteq G, X \neq \emptyset, X \notin \mathcal{F}_G\},$$

*that is,  $\mathcal{H}$  is a family of sub-itemsets of  $G$  not belonging to  $\mathcal{F}_G$ . Under the 0-hypothesis that the itemsets in  $\mathcal{H}$  are distributed according to  $p^*$ , the quantity  $2|D| r(G; \mathcal{F})$  is distributed asymptotically as  $\chi^2$  with degree  $|\mathcal{H}| = 2^{|G|} - 1 - |\mathcal{F}_G|$  of freedom.*

Theorem 5 is stated (but not proven) in a more general form in [23]. A rather technical proof is provided in Appendix A.

Theorem 5 motivates us to define the *normalised rank measure* to be the one-sided  $\chi^2$  test, that is,

$$nr(G; \mathcal{F}, D) = cdf(2|D| r(G; \mathcal{F}, D)),$$

where  $cdf(a) = P(\chi^2 < a)$  is the cumulative distribution function of  $\chi^2$  with degree  $2^{|G|} - 1 - |\mathcal{F}_G|$  of freedom. The number of degrees for different rank measures are provided in Table 1.

The following well-known result and its corollaries will play an important role in evaluating the measures.

**Lemma 6** *Let  $p^*$  be the maximum entropy distribution for itemsets  $\mathcal{F}$  and the corresponding frequencies  $\theta$ . Let  $q$  be a distribution satisfying the itemsets  $\mathcal{F}$ . Then we have*

$$-\sum_{\omega} q(\omega) \log p^*(\omega) = H(p^*).$$

**Corollary 7** *Let  $\mathcal{F}$  be the family of itemsets. We have that*

$$r(G; \mathcal{F}) = H(p^*) - H(q_G),$$

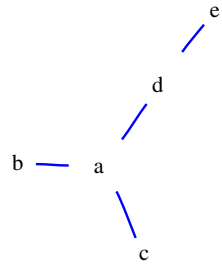
*where  $p^*$  is the maximum entropy distribution and  $q_G$  is the empirical distribution.*

**Corollary 8** *Let  $\mathcal{F}, \mathcal{H}$  be the families of itemsets such that  $\mathcal{H} \subseteq \mathcal{F}$ . Let  $p_1^*$  be the maximum entropy distribution for  $\mathcal{F}$  and let  $p_2^*$  be the maximum entropy distribution for  $\mathcal{H}$ . We have that*

$$r(G; \mathcal{F}) = \text{KL}(q_G \| p_2^*) - \text{KL}(p_1^* \| p_2^*),$$

*$q_G$  is the empirical distribution.*

**Fig. 1** A toy tree model. The related itemsets  $\{a, b, c, d, e, ab, ac, ad, de\}$  correspond to the attributes and the edges of the tree



**Corollary 9** Let  $\mathcal{F}, \mathcal{H}$  be the families of itemsets such that  $\mathcal{H} \subseteq \mathcal{F}$ . We have that

$$r(G; \mathcal{F}) \leq r(G; \mathcal{H}) .$$

### 3.3 Flexible models

So far we have considered ranks with fixed families of itemsets. In this section we introduce two additional models. In these models the itemsets are selected such that they minimise the rank.

Our first rank measure is the optimal tree model. A tree model can be described as a tree defined on the attributes of  $G$ . The corresponding family  $\mathcal{T}$  of itemsets contains the attributes from  $G$  and the itemsets of size 2 corresponding to the edges of the tree.

*Example 10* Consider  $G = \{a, b, c, d, e\}$  and consider the tree given in Fig. 1. The corresponding family of itemsets is  $\mathcal{T} = \{a, b, c, d, e, ab, ac, ad, de\}$ .

We can show that the Maximum Entropy distribution for  $\mathcal{T}$  has the form

$$p^* = \prod_{\{a,b\} \in \mathcal{T}} p^*(a, b) / \prod_{a \in G} p^*(a) .$$

This is, of course, Chow-Liu tree model [9]. We define the optimal tree to be the one that minimises the rank, that is,

$$\mathcal{T}^* = \underset{\mathcal{T} \text{ is a tree}}{\operatorname{arg\,min}} r(G; \mathcal{T}, D) .$$

To solve this tree let  $p_{ind}$  be the independence distribution. Corollary 8 allows us to rewrite the rank measure as

$$r(G; \mathcal{T}) = \operatorname{KL}(q_D \| p^*) = \operatorname{KL}(q_D \| p_{ind}) - \operatorname{KL}(p^* \| p_{ind}) .$$

Note that the first term  $\operatorname{KL}(q_D \| p_{ind})$  does not depend on  $\mathcal{T}$ . Hence we need to maximise the second term  $\operatorname{KL}(p^* \| p_{ind})$ . This is the mutual information of the tree and maximising this term is equivalent to finding maximum spanning tree in the mutual information graph. This can be done in polynomial time [9].

There is a deep connection between the rank  $r(G; \mathcal{T})$  and the rank for D-trees suggested in [19]. We can rewrite, by applying Corollary 7, the rank as

$$r(G; \mathcal{T}) = \operatorname{KL}(q_D \| p^*) = H(p^*) - H(q_G) .$$

The first term  $H(p^*)$  is the rank that is used in [19]. The authors in [19] seek patterns that have small  $H(p^*)$ , that is, trees that have strong dependencies between the attributes,

whereas we are interested in patterns that produce large  $r(G; T^*)$ , sets of attributes whose joint distribution cannot be explained even by the best tree model.

Our second model involves in finding a downward closed family  $\mathcal{F}$  of itemsets that produces the smallest normalised rank. Note that Corollary 9 implies that the rank decreases when we increase the number of known itemsets. However, this does not hold for the normalised rank and we will see that, contrary to the expectations, the best model can be different than  $\mathcal{A}_G$ , the set of all sub-itemsets of  $G$ . In other words, knowing all sub-itemsets does not guarantee the best model but, in fact, itemsets of higher order may mislead the prediction.

Unlike with the tree models, to our knowledge, there is no polynomial algorithm for finding the optimal downward closed family. Hence, we suggest a simple greedy approach. We start from the itemsets of size 1 and select the itemset from the negative border that minimises the rank. The itemset is added into the family and the procedure is repeated until there is no itemset that can decrease the rank. The algorithm is stated in Algorithm 1. We use  $\mathcal{F}^*$  to denote the resulting family.

---

**Algorithm 1** Greedy algorithm for finding the optimal downward closed family of item sets. The input is the data set  $D$  and the query itemset  $G$ . The output is  $\mathcal{F}^*$  a family of itemsets that produces low rank for the itemset

---

```

 $\mathcal{F}^* \leftarrow \mathcal{I}_G$ . {Initialise  $\mathcal{F}^*$  with itemsets of size 1.}
repeat
   $Y \leftarrow \arg \min_{X \in \text{negbord}(\mathcal{F}^*)} nr(G; \mathcal{F}^* \cup X)$ .
  if  $nr(G; \mathcal{F}^* \cup Y) < nr(G; \mathcal{F}^*)$  then
     $\mathcal{F}^* \leftarrow \mathcal{F}^* \cup Y$ .
  end if
until no more changes in  $\mathcal{F}^*$ .
    
```

---

### 3.4 Computing rank

Corollary 7 allows us to rewrite the rank as a difference of two entropies

$$r(G) = \text{KL}(q_G \| p^*) = H(p^*) - H(q_G).$$

Both distributions have  $|\Omega_G| = 2^{|G|}$  entries. However, the distribution  $q_G$  can have only  $|D|$  positive entries at maximum, hence the term  $H(q_G)$  can be computed efficiently.

The challenge in calculating the measure is to solve the Maximum Entropy distribution  $p^*$  and calculate its entropy. This can be done in polynomial time for the independence model and for the tree models. However, in the general case solving  $p^*$  is an **NP**-complete problem [10,30]; In such cases the distribution is solved using Iterative Scaling algorithm [12, 21]. The algorithm consists of consecutive steps. One such step requires  $O(|\Omega_G|) = O(2^{|G|})$  time. Hence computing the measure requires exponential time but it is doable for itemsets of reasonable size. The summary for evaluation times is provided in Table 1.

### 3.5 The effect of pruning itemsets

Note that in defining the measure we only use itemsets that are subsets of the query itemset  $G$ . This pruning guarantees that the number of entries in the distributions is  $2^{|G|}$  and not, at worst,  $2^K$ , where  $K$  is the number of columns in the dataset. Pruning attributes is essential



**Table 1** Summary of the rank measure

Measure	Description	# of degrees	Evaluation time
$r(G; \mathcal{I})$	Independence model	$2^{ G } - 1 -  G $	$O( G )$
$r(G; \mathcal{C})$	Gaussian model	$2^{ G } - 1 - \frac{1}{2} G  ( G  + 1)$	$O(2^{ G })$ per iter.
$r(G; \mathcal{A})$	All subsets model	1	$O(2^{ G })$ per iter.
$r(G; \mathcal{T}^*)$	Optimal tree model	$2^{ G } - 2 G $	$O( G ^2)$
$r(G; \mathcal{F}^*)$	Optimal family model	$2^{ G } - 1 -  \mathcal{F} $	$O(8^{ G })$ per iter.

The number of degrees, the third column, is used as a parameter for  $\chi^2$  distribution, when computing the normalised rank. The fourth column represents the evaluation times for the entropy of  $p^*$

since solving  $p^*$  is exponential to the number of attributes. The downside is that pruning may change the prediction as the following example demonstrates.

*Example 11* Assume that we have 3 attributes,  $a$ ,  $b$ , and  $c$ . Our known itemsets are  $\mathcal{F} = \{a, b, c, ac, bc\}$  and their frequencies are  $\theta_a = \theta_b = \theta_c = \theta_{ac} = \theta_{bc} = 1/2$ . In other words, the attributes are identical and correspond to a fair coin flip. Assume that we are interested in rank of  $G = ab$ . In this case the pruned family of itemsets is  $\mathcal{F}_G = \{a, b\}$  and the Maximum Entropy distribution is the uniform distribution. The empirical distribution is

$$q_G(a = 0, b = 0) = q_G(a = 1, b = 1) = 1/2$$

$$q_G(a = 1, b = 0) = q_G(a = 0, b = 1) = 0.$$

The rank is then  $r(ab; \mathcal{F}) = 0.69$ . However, if we had used the frequencies of  $ac$  and  $bc$ , we would have concluded that  $a = b$  and that the Maximum Entropy distribution is equal to the empirical distribution, hence the rank would have been 0.

In [31] we investigate the effect of pruning attributes and conclude that in some cases we can remove a large portion of attributes outside  $G$ . However, in those cases, the family of known itemsets has many restrictions and, for instance, we cannot remove safely any attribute from the Gaussian model.

### 4 Related work

Traditionally, the support (frequency) of the itemset is used for ranking itemsets. Alternative measures that resemble the support are studied in [26].

Our work resembles approach of [6] in which the authors defined the significance of an itemset by comparing the distribution  $q_G$  against the independence model. The authors used  $\chi^2$  statistical test as a measure, that is, if  $p$  is the distribution related to the independence model, the rank measure is

$$r_b(G) = \sum_{\omega \in \Omega_G} \frac{(q_G(\omega) - p(\omega))^2}{p(\omega)}. \tag{1}$$

In [14] the authors also compare the frequency of an itemset against the independence model but in addition they use Bayes screening to smooth the values. Also, in [1] the authors proposed the collective strength as a measure of significance. To be more specific, we say

that a transaction  $\omega \in \Omega_G$  is *good* if it contains only 0s or only 1s. Let  $p$  be the distribution related to the independence model. Then the measure is

$$r_{cs}(G) = \frac{q_G(\omega \text{ is good})}{p(\omega \text{ is good})} \frac{p(\omega \text{ is bad})}{q_G(\omega \text{ is bad})}. \quad (2)$$

This measure obtains small values when data obeys the independence model. In a related work presented in [13] the authors define an itemset to be interesting if its frequency increases significantly from one dataset to another. In [17] the authors order itemsets based on their p-values. In [19] the authors used entropy of tree models for ranking itemsets. In addition, many measures has been suggested for ranking association rules [2, 7, 20, 29].

The authors in [28] showed empirically that Maximum entropy model provides excellent estimates for itemsets. Rank can be used for pruning a large family of itemsets by picking the itemsets having the largest rank. Other pruning methods are proposed in [4, 8, 27]. The authors in [34] suggest a generic framework for discovering significant rules. In addition, a relevant framework is described in [24]; the authors define a pattern ordering given an estimation algorithm and a loss function. In [25] the authors use information component analysis to find patterns in a drug safety database.

## 5 Experiments

In this section we present our empirical results. In the first three sections we explain the datasets and the setup. In our experiments we investigate the significance of itemsets, how different measures are related to each other, and the monotonicity of the ranks.

### 5.1 Synthetic datasets

For the testing purposes we created two synthetic datasets. Each dataset contained 100 attributes and 5,000 rows. The first dataset, *gen-ind*, was generated such that the attributes were independent. The margins were sampled uniformly from [0, 1]. In the second dataset, *gen-copy*, each column was a copy of the previous column corrupted by the symmetric white noise. The amount of noise, that is the probability

$$p(a_i = 1 \mid a_{i-1} = 0) = p(a_i = 0 \mid a_{i-1} = 1),$$

was selected uniformly from [0, 1] for each column  $a_i$ , individually. The first column was generated by a coin flip. Our expectations are that in *gen-ind* the itemsets of size 1 are significant and that in *gen-copy* the itemsets of size 2 are significant.

### 5.2 Real datasets

In our experiments we used the following real-world datasets. Data in *Accidents*<sup>1</sup> were obtained from the Belgian “Analysis Form for Traffic Accidents” forms that is filled out by a police officer for each traffic accident that occurs with injured or deadly wounded casualties on a public road in Belgium. In total, 340,183 traffic accident records are included in the dataset [18]. The datasets *POS*<sup>2</sup>, *WebView-1*<sup>3</sup> and *WebView-2*<sup>4</sup> were contributed by

<sup>1</sup> <http://fimi.cs.helsinki.fi/data/accidents.dat.gz>

<sup>2</sup> <http://www.ecn.purdue.edu/KDDCUP/data/BMS-POS.dat.gz>

<sup>3</sup> <http://www.ecn.purdue.edu/KDDCUP/data/BMS-WebView-1.dat.gz>

<sup>4</sup> <http://www.ecn.purdue.edu/KDDCUP/data/BMS-WebView-2.dat.gz>

Blue Martini Software as the KDD Cup 2000 data [22]. *POS* contains several years worth of point-of-sale data from a large electronics retailer. *WebView-1* and *WebView-2* contain several months worth of click-stream data from two e-commerce web sites. *Kosarak*<sup>5</sup> consists of (anonymised) click-stream data of a Hungarian on-line news portal. *Retail*<sup>6</sup> is a retail market basket data supplied by an anonymous Belgian retail supermarket store [5]. The dataset *Paleo*<sup>7</sup> contains information of species fossils found in specific paleontological sites in Europe [15], preprocessed as in [16].

### 5.3 Setup for the experiments

In this section we will describe how we conducted our experiments. We reduced the largest datasets by selecting the first 10,000 rows and 200 most frequent attributes. From each dataset we computed all *almost non-derivable* itemsets. By almost non-derivable we mean that the difference between the upper bound and the lower bound of a given itemset, say  $G$ , is at least  $n$  transactions. In other words, if we know the frequencies of all sub-itemsets of  $G$ , then we cannot predict the frequency of  $G$  within  $n$  transactions. If  $n = 0$ , then an itemset is non-derivable. It is known that the family of almost non-derivable itemsets is anti-monotonic [8, Lemma 3.1]. A reason to use almost non-derivable itemsets instead of frequent itemsets is the statement of Theorem 3, that is,  $r(G; \mathcal{A}) = 0$  if the itemset is derivable. The other reason is that we want to study how the measure behaves for infrequent itemsets.

To keep the sizes of the obtained families within reasonable bounds we used different thresholds for different datasets: For *gen-ind*, *Retail* and *WebView-2* we set  $n = 5$ . For *POS* the threshold  $n$  was set to 10 and for *gen-copy* and *Accidents*  $n$  was set to 100. For the rest of the datasets we set  $n = 0$ , that is, we mined all non-derivable itemsets from these datasets.

For each itemset from the obtained itemsets we queried the following measures:

- Frequency.
- Normalised rank measures  $nr(G; \mathcal{I})$ ,  $nr(G; \mathcal{C})$ ,  $nr(G; \mathcal{A})$ ,  $nr(G; \mathcal{T}^*)$ ,  $nr(G; \mathcal{F}^*)$ .
- Measures discussed in Sect. 4: A  $\chi^2$  test  $r_b(G)$  defined in Eq. 1 and a collective strength  $r_{cs}(G)$  defined in Eq. 2.

The evaluation times and the sizes of the query families are given in Table 2.

### 5.4 Significant itemsets

Our first experiment is to study how many of the itemsets are significant. We did this by comparing our rank measures with risk level 0.05. The results are given in Tables 3, 4 and 5. We also provide a typical example of box plots in Fig. 2.

Let us first study *gen-ind*, a synthetic dataset with independent columns. We see from Table 3 that according to  $nr(G; \mathcal{I})$  a large portion of itemsets of size 1 are significant but only a small portion of itemsets having size larger than 1 is significant. This is an expected result since the frequencies obey the independence model. In Tables 4 we have similar results for  $nr(G; \mathcal{C})$  and for  $nr(G; \mathcal{A})$ . However, the values of  $nr(G; \mathcal{C})$  and for  $nr(G; \mathcal{A})$  tend to be larger than the values of  $nr(G; \mathcal{I})$ . The reason for this is a type of overlearning: Since the frequencies of itemsets are calculated from the datasets, they are imprecise. Hence, the itemsets of larger size mislead us during prediction, because the resulting Maximum Entropy distribution is not an independent model (although close to one).

<sup>5</sup> <http://fimi.cs.helsinki.fi/data/kosarak.dat.gz>

<sup>6</sup> <http://fimi.cs.helsinki.fi/data/retail.dat.gz>

<sup>7</sup> NOW public release 030717 available from [15].

**Table 2** The evaluation times and the sizes of the query families

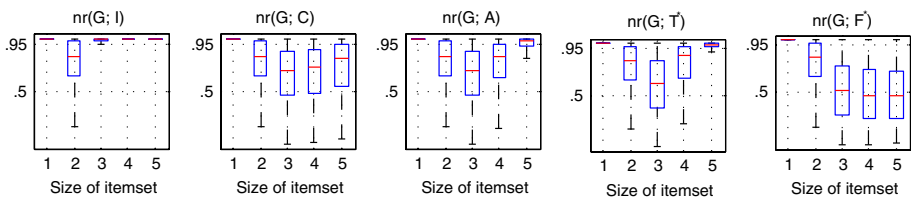
Data	$n$	# of $G$	$\max G $	Evaluation times				
				$nr(G; \mathcal{I})$ (s)	$nr(G; \mathcal{C})$	$nr(G; \mathcal{A})$	$nr(G; \mathcal{T}^*)$ (s)	$nr(G; \mathcal{F}^*)$
<i>gen-ind</i>	5	156,699	6	2	52 s	29 min	2	11 min
<i>gen-copy</i>	100	111,487	4	0	12 s	57 s	0	1 min
<i>Accidents</i>	100	354,399	6	2	1 min	19 min	3	108 min
<i>Kosarak</i>	5	223,734	5	1	4 s	9 s	0	47 s
<i>Paleo</i>	0	166,903	5	0	8 s	35 s	0	2 min
<i>POS</i>	10	246,640	6	1	8 s	27 s	1	5 min
<i>Retail</i>	0	818,813	6	3	19 s	49 s	4	4 min
<i>WebView-1</i>	5	226,313	5	1	5 s	8 s	1	39 s
<i>WebView-2</i>	0	715,398	6	3	27 s	2 min	4	11 min

The second column is the threshold used in mining almost non-derivable itemsets. The fourth column is the maximal size of a query itemset. The evaluation time does not include the time spent mining the itemsets

**Table 3** The percentages of significant itemsets according to  $nr(G; \mathcal{I})$

Data	itemset size							All
	1	2	3	4	5	6		
<i>gen-ind</i>	0.92	0.05	0.04	0.03	0.02	0.01	0.03	
<i>gen-copy</i>	0.08	0.14	0.24	0.03	–	–	0.07	
<i>Accidents</i>	0.99	0.60	0.95	1	1	1	0.97	
<i>Kosarak</i>	1	0.62	0.99	1	1	–	0.96	
<i>Paleo</i>	1	0.30	0.81	0.99	1	–	0.88	
<i>POS</i>	1	0.45	0.99	1	1	1	0.95	
<i>Retail</i>	1	0.14	0.30	0.93	1	1	0.45	
<i>WebView-1</i>	1	0.70	1	1	1	–	0.97	
<i>WebView-2</i>	1	0.20	0.69	1	1	1	0.85	

Each entry is a fraction of itemsets of specific size calculated from a specific dataset. Significance is measured using  $\chi^2$  distribution with 0.05 risk level



**Fig. 2** Box plots of the rank measures computed from *Paleo*

Let us continue by studying *gen-copy*, a synthetic data in which an attribute is a noisy copy of the previous attribute. We see that  $nr(G; \mathcal{T}^*)$  tends to have smaller ranks than  $nr(G; \mathcal{I})$  when  $G$  has size 3. The reason for this is that, unlike with *gen-ind*, the independence model cannot explain the dataset. However, when we predict using also the itemsets of size 2, the

**Table 4** The percentages of significant itemsets according to  $nr(G; C)$  and  $nr(G; \mathcal{A})$

Data	$nr(G; C)$ , itemset size							$nr(G; \mathcal{A})$ , itemset size						
	1	2	3	4	5	6	All	1	2	3	4	5	6	All
<i>gen-ind</i>	0.92	0.05	0.06	0.05	0.04	0.03	0.05	0.92	0.05	0.06	0.06	0.06	0.06	0.06
<i>gen-copy</i>	0.08	0.14	0.06	0.03	–	–	0.03	0.08	0.14	0.06	0.05	–	–	0.05
<i>Accidents</i>	0.99	0.60	0.21	0.45	0.62	0.60	0.45	0.99	0.60	0.21	0.07	0.05	0.06	0.11
<i>Kosarak</i>	1	0.62	0.32	0.50	0.38	–	0.37	1	0.62	0.32	0.10	0.04	–	0.33
<i>Paleo</i>	1	0.30	0.12	0.15	0.21	–	0.15	1	0.30	0.12	0.21	0.64	–	0.18
<i>POS</i>	1	0.45	0.09	0.21	0.43	0.66	0.17	1	0.45	0.09	0.06	0.07	0.05	0.11
<i>Retail</i>	1	0.14	0.04	0.08	0.12	0.38	0.05	1	0.14	0.04	0.15	0.27	0.25	0.07
<i>WebView-1</i>	1	0.70	0.48	0.32	0.52	–	0.48	1	0.70	0.48	0.09	0.07	–	0.45
<i>WebView-2</i>	1	0.20	0.11	0.20	0.88	1	0.17	1	0.20	0.11	0.16	0.36	0.48	0.15

Each entry is a fraction of itemsets of specific size calculated from a specific dataset. Significance is measured using  $\chi^2$  distribution with 0.05 risk level

**Table 5** The percentages of significant itemsets according to  $nr(G; T^*)$  and  $nr(G; \mathcal{F}^*)$

Data	$nr(G; T^*)$ , itemset size							$nr(G; \mathcal{F}^*)$ , itemset size						
	1	2	3	4	5	6	All	1	2	3	4	5	6	All
<i>gen-ind</i>	0.92	0.05	0.02	0.01	0.01	0	0.01	0.92	0.05	0.01	0.01	0	0	0.01
<i>gen-copy</i>	0.08	0.14	0.02	0	–	–	0.01	0.08	0.14	0.01	0	–	–	0.01
<i>Accidents</i>	0.99	0.60	0.40	0.80	0.95	0.97	0.75	0.99	0.60	0.18	0.12	0.13	0.02	0.15
<i>Kosarak</i>	1	0.62	0.80	0.94	1	–	0.80	1	0.62	0.32	0.05	0.03	–	0.32
<i>Paleo</i>	1	0.30	0.10	0.35	0.81	–	0.24	1	0.30	0.06	0.05	0.04	–	0.07
<i>POS</i>	1	0.45	0.47	0.99	1	1	0.65	1	0.45	0.08	0.02	0.02	0	0.09
<i>Retail</i>	1	0.14	0.03	0.18	0.78	1	0.07	1	0.14	0.01	0.02	0.02	0.13	0.02
<i>WebView-1</i>	1	0.70	0.83	1	1	–	0.84	1	0.70	0.46	0.07	0.30	–	0.43
<i>WebView-2</i>	1	0.20	0.11	0.57	1	1	0.37	1	0.20	0.06	0.03	0.14	0.44	0.05

Each entry is a fraction of itemsets of specific size calculated from a specific dataset. Significance is measured using  $\chi^2$  distribution with 0.05 risk level

prediction becomes more accurate. The measures  $nr(G; C)$  and  $nr(G; \mathcal{A})$  also produce small ranks, however, these ranks tend to be slightly larger than the ranks of  $nr(G; T^*)$ .

We turn our attention to real datasets. We see that for these datasets the independence model is too strict: According to  $nr(G; \mathcal{T})$  almost all itemsets are significant: The results change drastically, when we use richer models. According to  $nr(G; \mathcal{A})$  only about 5–50% of the itemsets are significant, depending on the dataset. Similar overfitting that occurred with *gen-ind* also occurs in some but not all real datasets (see Fig. 2). For instance, in *Retail*  $nr(G; \mathcal{A})$  tends to produce higher values than  $nr(G; C)$  but not in *POS*.

### 5.5 The effect of the known itemsets

We continued our experiments by comparing the measures  $nr(G; \mathcal{T})$ ,  $nr(G; C)$ ,  $nr(G; \mathcal{A})$ ,  $nr(G; T^*)$ , and  $nr(G; \mathcal{F}^*)$  against each other. This was done by calculating the correlations between the rank measures. The results are given in Tables 6 and 7.

**Table 6** Correlations between the measures  $nr(G; \mathcal{I})$ ,  $nr(G; \mathcal{C})$ , and  $nr(G; \mathcal{A})$

Data	$nr(G; \mathcal{I})$ vs. $nr(G; \mathcal{C})$	$nr(G; \mathcal{I})$ vs. $nr(G; \mathcal{A})$	$nr(G; \mathcal{C})$ vs. $nr(G; \mathcal{A})$
<i>gen-ind</i>	0.74	0.26	0.39
<i>gen-copy</i>	0.52	0.28	0.53
<i>Accidents</i>	0.17	0.07	0.37
<i>Kosarak</i>	0.14	0.13	0.90
<i>Paleo</i>	0.16	0.22	0.67
<i>POS</i>	0.12	0.10	0.77
<i>Retail</i>	0.62	0.67	0.88
<i>WebView-1</i>	0.14	0.12	0.88
<i>WebView-2</i>	0.43	0.51	0.68

**Table 7** Correlations between the flexible measures  $nr(G; T^*)$  and  $nr(G; \mathcal{F}^*)$  and the measures  $nr(G; \mathcal{I})$ ,  $nr(G; \mathcal{C})$ , and  $nr(G; \mathcal{A})$

Data	$nr(G; T^*)$ vs.			$nr(G; \mathcal{F}^*)$ vs.			
	$nr(G; \mathcal{I})$	$nr(G; \mathcal{C})$	$nr(G; \mathcal{A})$	$nr(G; \mathcal{I})$	$nr(G; \mathcal{C})$	$nr(G; \mathcal{A})$	$nr(G; T^*)$
<i>gen-ind</i>	0.81	0.94	0.37	0.81	0.91	0.36	0.98
<i>gen-copy</i>	0.62	0.92	0.49	0.64	0.89	0.49	0.97
<i>Accidents</i>	0.34	0.57	0.20	0.07	0.58	0.54	0.29
<i>Kosarak</i>	0.41	0.41	0.36	0.13	0.85	0.93	0.46
<i>Paleo</i>	0.35	0.64	0.47	0.16	0.84	0.67	0.62
<i>POS</i>	0.42	0.34	0.28	0.05	0.72	0.84	0.21
<i>Retail</i>	0.66	0.82	0.82	0.59	0.92	0.84	0.86
<i>WebView-1</i>	0.56	0.29	0.24	0.12	0.86	0.93	0.28
<i>WebView-2</i>	0.59	0.62	0.65	0.34	0.77	0.72	0.54

From the results we see that all correlations are positive. For the real datasets the correlations between  $nr(G; \mathcal{C})$  and  $nr(G; \mathcal{A})$  are systematically higher than the correlations between  $nr(G; \mathcal{I})$  and  $nr(G; \mathcal{A})$  or between  $nr(G; \mathcal{C})$  and  $nr(G; \mathcal{A})$ . This implies that  $nr(G; \mathcal{I})$  produces different ranks whereas  $nr(G; \mathcal{C})$  and  $nr(G; \mathcal{A})$  are more similar. This supports the behaviour we have seen in Sect. 5.4.

The measure  $nr(G; \mathcal{F}^*)$  correlate more with  $nr(G; \mathcal{A})$  and  $nr(G; \mathcal{C})$  than with  $nr(G; \mathcal{I})$ . The correlation between  $nr(G; \mathcal{F}^*)$  and  $nr(G; T^*)$  is somewhat weaker but it is stronger than the correlation between  $nr(G; \mathcal{F}^*)$  and  $nr(G; \mathcal{I})$ .

### 5.6 Flexible models

Our next goal is to compare the flexible measures  $nr(G; T^*)$  and  $nr(G; \mathcal{F}^*)$  against the rest of the measures. From Table 5 we see that  $nr(G; \mathcal{F}^*)$  tend to produce the smallest amount of significant itemsets whereas the  $nr(G; T^*)$  produces large ranks, especially for queries with many attributes.

We calculated the number of queries in which  $nr(G; T^*)$  and  $nr(G; \mathcal{F}^*)$  produce smaller rank than the rest of the measures. Since the measures are equivalent for the queries of size 1

**Table 8** Percentages of queries in which the flexible measures  $nr(G; \mathcal{T}^*)$  and  $nr(G; \mathcal{F}^*)$  outperform the other rank measures

Data	$nr(G; \mathcal{T}^*) \leq$			$nr(G; \mathcal{F}^*) \leq$			
	$nr(G; \mathcal{I})$	$nr(G; \mathcal{C})$	$nr(G; \mathcal{A})$	$nr(G; \mathcal{I})$	$nr(G; \mathcal{C})$	$nr(G; \mathcal{A})$	$nr(G; \mathcal{T}^*)$
<i>gen-ind</i>	0.84	0.97	0.78	1	0.99	0.81	0.95
<i>gen-copy</i>	0.79	0.96	0.82	1	0.99	0.86	0.99
<i>Accidents</i>	1	0.16	0.13	1	0.94	0.66	0.95
<i>Kosarak</i>	1	0.08	0.08	1	0.38	0.36	0.96
<i>Paleo</i>	0.99	0.47	0.58	1	0.91	0.85	0.86
<i>POS</i>	1	0.12	0.12	1	0.65	0.62	0.90
<i>Retail</i>	0.92	0.77	0.80	1	0.94	0.92	0.62
<i>WebView-1</i>	1	0.19	0.19	1	0.52	0.49	0.93
<i>WebView-2</i>	0.99	0.37	0.46	1	0.89	0.87	0.79

Queries only of size 3 or larger were considered

**Table 9** Number of itemsets occurring in  $\mathcal{F}^*$ , the family of known itemsets in  $r(G; \mathcal{F}^*)$ , normalised by the maximum number of possible occurrences

Data	Ratio of used itemsets			
	2	3	4	5
<i>gen-ind</i>	0.35	0.01	0	0
<i>gen-copy</i>	0.36	0.01	–	–
<i>Accidents</i>	0.87	0.36	0.02	0
<i>Kosarak</i>	0.95	0.49	0.01	–
<i>Paleo</i>	0.74	0.16	0	–
<i>POS</i>	0.96	0.47	0.01	0
<i>Retail</i>	0.62	0.13	0	0
<i>WebView-1</i>	0.93	0.44	0	–
<i>WebView-2</i>	0.81	0.26	0.07	0

Each column represent itemsets of specific size

and 2, these queries were ignored. From the results given in Table 8 we see that the flexible models outperform  $nr(G; \mathcal{I})$ , however, the performance against other measure depends on the data set. For instance,  $nr(G; \mathcal{F}^*)$  outperform  $nr(G; \mathcal{C})$  and  $nr(G; \mathcal{A})$  in *Retail* but produces larger ranks in *Kosarak*. This suggests that the greedy algorithm sometimes fails to find the optimal family  $\mathcal{F}^*$ .

We studied the sizes of itemsets occurring in  $\mathcal{F}^*$ , the family of known itemsets in  $nr(G; \mathcal{F}^*)$ . To be more precise, let  $\mathcal{F}_G^*$  be the family of known itemsets for the query  $G$ . Let  $L$  be the size of itemsets we are interested in. We define the ratio  $r_L$  to be

$$r_L = \frac{\sum_G |\{X \in \mathcal{F}_G^*; |X| = L\}|}{\sum_G \binom{|G|}{L}},$$

that is, the number of itemset of size  $L$  occurring in  $\mathcal{F}^*$  divided by the maximum number of occurrences. The ratios  $r_L$  are given in Table 9. We see that the itemsets of size 2 and 3 are frequently used, however, the itemsets of larger size are rarely used.

**Table 10** Correlations between the rank measures  $nr(G; \mathcal{I})$ ,  $nr(G; \mathcal{C})$ , and  $nr(G; \mathcal{A})$  and the base measures: the frequency of  $G$ ,  $r_b(G)$ , the  $\chi^2$  test for independency, and  $r_{cs}(G)$ , the collective strength of the itemset  $G$

Data	$nr(G; \mathcal{I})$ vs.			$nr(G; \mathcal{C})$ vs.			$nr(G; \mathcal{A})$ vs.		
	Freq.	$r_b(G)$	$r_{cs}(G)$	Freq.	$r_b(G)$	$r_{cs}(G)$	Freq.	$r_b(G)$	$r_{cs}(G)$
<i>gen-ind</i>	0.06	0.99	-0.01	0.03	0.72	-0.01	0	0.25	0
<i>gen-copy</i>	0.15	1	0.02	0.07	0.52	0.02	0.01	0.27	0.01
<i>Accidents</i>	0.01	1	0.02	-0.01	0.17	0.05	0.03	0.07	0.01
<i>Kosarak</i>	0.01	0.98	0.20	0.01	0.14	0.27	0	0.13	0.21
<i>Paleo</i>	0.18	0.95	0.39	0.01	0.15	0.10	-0.03	0.20	0.03
<i>POS</i>	0.05	0.99	0.22	0.09	0.12	0.20	0.07	0.10	0
<i>Retail</i>	0.04	0.97	0.31	0.05	0.57	0.17	0.05	0.61	0.25
<i>WebView-1</i>	0.06	0.98	0.19	0.07	0.15	-0.29	0.05	0.13	-0.32
<i>WebView-2</i>	0.12	0.96	0.33	0.17	0.36	0.39	0.12	0.43	0.25

**Table 11** Correlations between the rank measures  $nr(G; \mathcal{T}^*)$  and  $nr(G; \mathcal{F}^*)$  and the base measures: the frequency of  $G$ ,  $r_b(G)$ , the  $\chi^2$  test for independency, and  $r_{cs}(G)$ , the collective strength of the itemset  $G$

Data	$nr(G; \mathcal{T}^*)$ vs.			$nr(G; \mathcal{F}^*)$ vs.		
	Freq.	$r_b(G)$	$r_{cs}(G)$	Freq.	$r_b(G)$	$r_{cs}(G)$
<i>gen-ind</i>	0.07	0.79	-0.01	0.06	0.79	-0.01
<i>gen-copy</i>	0.16	0.62	0.03	0.16	0.64	0.03
<i>Accidents</i>	-0.02	0.33	0.04	0.04	0.07	0.01
<i>Kosarak</i>	0.01	0.39	0.32	0	0.13	0.25
<i>Paleo</i>	0.24	0.28	0.44	0.09	0.12	0.03
<i>POS</i>	0.12	0.41	0.37	0.07	0.05	-0.23
<i>Retail</i>	0.06	0.58	0.31	0.06	0.53	0.14
<i>WebView-1</i>	0.11	0.55	0.16	0.04	0.13	-0.35
<i>WebView-2</i>	0.20	0.49	0.46	0.16	0.28	0.14

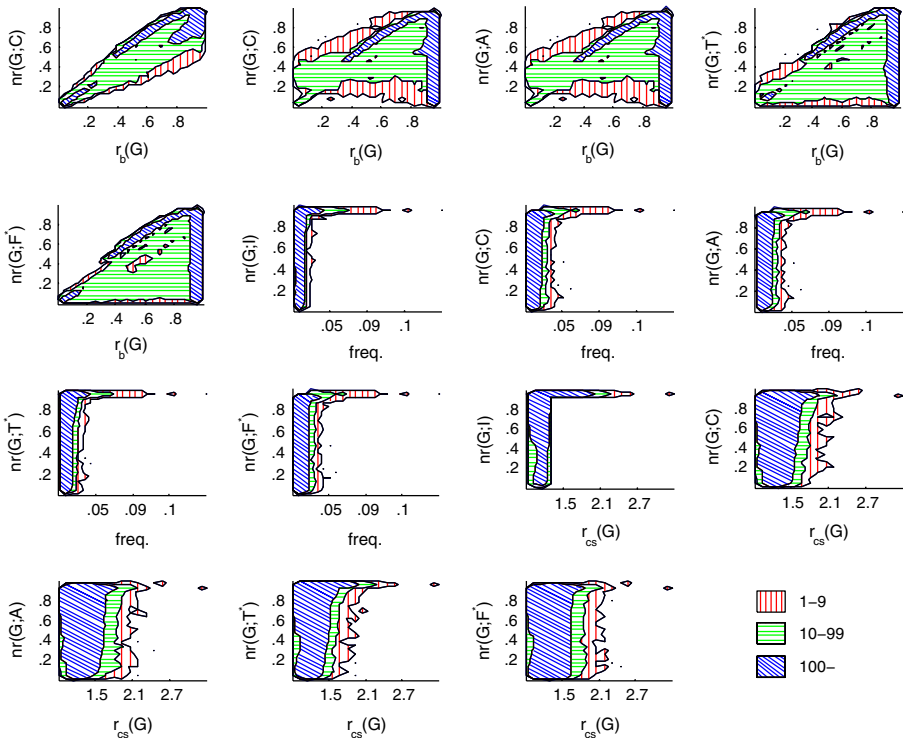
### 5.7 Rank versus other methods

We compared our measures against the other ranking methods described in Sect. 5.3. Namely, we calculated the correlations of  $nr(G; \mathcal{I})$ ,  $nr(G; \mathcal{C})$ ,  $nr(G; \mathcal{A})$ ,  $nr(G; \mathcal{T}^*)$ , and  $nr(G; \mathcal{F}^*)$  against the frequency of  $G$ ,  $r_b(G)$ , the  $\chi^2$  test for independency, and  $r_{cs}(G)$ , the collective strength of the itemset  $G$ . The results are presented in Tables 10 and 11. We also studied the relationships by plotting our measures as functions of the aforementioned approaches and such examples are given in Fig. 3.

Our first observation is that  $nr(G; \mathcal{I})$  correlates strongly with  $r_b(G)$ . This is an expected result since both test the independency of attributes inside the itemsets and also because  $nr(G; \mathcal{I})$  is asymptotically a  $\chi^2$  test (see Theorem 5). There is some correlation between  $r_b(G)$  and the rest of the measures although this correlation is much weaker compared to  $nr(G; \mathcal{I})$ .

Apart from *WebView-2*, there is little correlation between the measures and the frequency.





**Fig. 3** Ranks as functions of the base measures. The plots are calculated from *Paleo* dataset

The correlation between the measures and the collective strength  $r_{cs}(G)$  exists but varies depending on the method and the dataset. The strongest correlations are obtained when  $r_{cs}(G)$  is compared against  $nr(G; \mathcal{I})$  or  $nr(G; T^*)$ . The dependency between  $nr(G; \mathcal{I})$  and  $r_{cs}(G)$  is a natural result since  $r_{cs}(G)$  produces small values when attributes are independent.

### 5.8 Monotonicity of rank

In this section we investigate the relationship between the rank of an itemset and the ranks of its sub-itemsets. Namely, we tested whether the measures are monotonic, that is, whether  $nr(G; \mathcal{F}) \geq nr(H; \mathcal{F})$  for all  $H \subset G$ . We deliberately ignored sub-itemsets having size 1 since they all have very high rank. We also tested whether the measures are anti-monotonic, that is, decreasing w.r.t. set inclusion.

From the results given in Tables 12, 13, 14 and 15 our first observation is that  $nr(G; \mathcal{I})$  are increasing for real datasets but not for the synthetic datasets. The raw values of  $nr(G; \mathcal{I})$  are indeed increasing but this does not hold for the  $P$ -values since the number of degrees varies. The measure  $nr(G; T^*)$  tends also be monotonic but not as much as  $nr(G; \mathcal{I})$ . On the contrary,  $nr(G; \mathcal{C})$ ,  $nr(G; \mathcal{A})$ , and  $nr(G; \mathcal{F}^*)$  are increasing for extremely few itemsets.

Table 14 suggests that  $nr(G; \mathcal{C})$ ,  $nr(G; \mathcal{A})$ , and  $nr(G; \mathcal{F}^*)$  satisfies the anti-monotonicity to some degree. Measures  $nr(G; \mathcal{C})$  and  $nr(G; \mathcal{A})$  are anti-monotonic for relatively high percentage of itemsets of size 3. Among itemsets of size 4,  $nr(G; \mathcal{F}^*)$  satisfies the property

**Table 12** Percentages of itemsets satisfying the property of monotonicity

Data	$nr(G; \mathcal{I})$					$nr(G; \mathcal{C})$					$nr(G; \mathcal{A})$				
	3	4	5	6	All	3	4	5	6	All	3	4	5	6	All
<i>gen-ind</i>	0.09	0.02	0.01	0	0.03	0.27	0.03	0.01	0	0.05	0.27	0.11	0.06	0.03	0.11
<i>gen-copy</i>	0.15	0.01	-	-	0.04	0.20	0.02	-	-	0.05	0.20	0.10	-	-	0.12
<i>Accidents</i>	0.78	0.92	0.97	0.99	0.90	0.01	0.02	0.02	0	0.02	0.01	0	0	0	0
<i>Kosarak</i>	0.93	0.98	1	-	0.93	0	0	0	-	0	0	0	0	-	0
<i>Paleo</i>	0.40	0.61	0.84	-	0.51	0.04	0	0	-	0.02	0.04	0	0	-	0.02
<i>POS</i>	0.87	1	1	1	0.92	0	0	0.01	0	0	0	0	0	0	0
<i>Retail</i>	0.11	0.42	0.92	1	0.19	0.04	0	0	0	0.03	0.04	0.02	0	0	0.04
<i>WebView-1</i>	0.98	1	1	-	0.98	0.04	0	0	-	0.04	0.04	0	0	-	0.04
<i>WebView-2</i>	0.39	0.88	1	1	0.67	0.04	0	0.08	1	0.02	0.04	0	0	0	0.02

The itemset  $G$  satisfies the property if  $nr(G; \mathcal{F}) \geq nr(H; \mathcal{F})$  for all  $H \subset G$  such that  $|H| \geq 2$

**Table 13** Percentages of itemsets satisfying the property of monotonicity

Data	$nr(G; \mathcal{T}^*)$					$nr(G; \mathcal{F}^*)$				
	3	4	5	6	All	3	4	5	6	All
<i>gen-ind</i>	0.13	0.01	0	0	0.02	0.07	0.01	0	0	0.02
<i>gen-copy</i>	0.10	0.01	-	-	0.02	0.05	0.01	-	-	0.02
<i>Accidents</i>	0.02	0.13	0.35	0.47	0.17	0.01	0	0	0	0
<i>Kosarak</i>	0	0.26	0.32	-	0.03	0	0	0	-	0
<i>Paleo</i>	0.02	0	0	-	0.01	0.02	0	0	-	0.01
<i>POS</i>	0	0.23	0.97	1	0.11	0	0	0	0	0
<i>Retail</i>	0.03	0	0.14	1	0.02	0.02	0	0	0	0.02
<i>WebView-1</i>	0.04	0.07	0.89	-	0.05	0.03	0	0	-	0.03
<i>WebView-2</i>	0.03	0.02	0.72	1	0.04	0.03	0	0	0	0.01

The itemset  $G$  satisfies the property if  $nr(G; \mathcal{F}) \geq nr(H; \mathcal{F})$  for all  $H \subset G$  such that  $|H| \geq 2$

**Table 14** Percentages of itemsets satisfying the property of anti-monotonicity

Data	$nr(G; \mathcal{I})$					$nr(G; \mathcal{C})$					$nr(G; \mathcal{A})$				
	3	4	5	6	All	3	4	5	6	All	3	4	5	6	All
<i>gen-ind</i>	0.21	0.07	0.03	0.02	0.07	0.25	0.07	0.03	0.01	0.07	0.25	0.08	0.02	0.01	0.07
<i>gen-copy</i>	0.15	0.06	-	-	0.08	0.25	0.08	-	-	0.11	0.25	0.07	-	-	0.10
<i>Accidents</i>	0.03	0	0	0	0.01	0.62	0.04	0	0	0.16	0.62	0.21	0.05	0.02	0.26
<i>Kosarak</i>	0.02	0.06	0	-	0.02	0.93	0.03	0.01	-	0.83	0.93	0.33	0.08	-	0.86
<i>Paleo</i>	0.02	0	0	-	0.01	0.43	0.04	0	-	0.22	0.43	0.07	0	-	0.23
<i>POS</i>	0.01	0.01	0.04	0.23	0.01	0.87	0.07	0	0	0.57	0.87	0.18	0.06	0.09	0.61
<i>Retail</i>	0.17	0	0	0	0.13	0.38	0.05	0.01	0	0.30	0.38	0.03	0.01	0	0.29
<i>WebView-1</i>	0	0	0	-	0	0.69	0.11	0	-	0.62	0.69	0.39	0.15	-	0.66
<i>WebView-2</i>	0.07	0.01	0.14	0.96	0.04	0.48	0.06	0	0	0.24	0.48	0.06	0.07	0.04	0.25

The itemset  $G$  satisfies the property if  $nr(G; \mathcal{F}) \leq nr(H; \mathcal{F})$  for all  $H \subset G$  such that  $|H| \geq 2$

**Table 15** Percentages of itemsets satisfying the property of anti-monotonicity

Data	$nr(G; \mathcal{T}^*)$					$nr(G; \mathcal{F}^*)$				
	3	4	5	6	All	3	4	5	6	All
<i>gen-ind</i>	0.45	0.16	0.06	0.02	0.15	0.53	0.16	0.05	0.01	0.15
<i>gen-copy</i>	0.43	0.18	–	–	0.22	0.51	0.17	–	–	0.23
<i>Accidents</i>	0.58	0.01	0	0	0.13	0.69	0.29	0.07	0.02	0.32
<i>Kosarak</i>	0.91	0	0	–	0.81	0.95	0.50	0.07	–	0.90
<i>Paleo</i>	0.52	0.02	0	–	0.25	0.56	0.15	0.01	–	0.34
<i>POS</i>	0.88	0	0	0	0.56	0.90	0.47	0.13	0	0.73
<i>Retail</i>	0.72	0.02	0	0	0.55	0.75	0.13	0.03	0	0.60
<i>WebView-1</i>	0.62	0	0	–	0.55	0.70	0.56	0.22	–	0.68
<i>WebView-2</i>	0.69	0.01	0	0	0.31	0.71	0.23	0.06	0	0.44

The itemset  $G$  satisfies the property if  $nr(G; \mathcal{F}) \leq nr(H; \mathcal{F})$  for all  $H \subset G$  such that  $|H| \geq 2$

of anti-monotonicity for a slightly larger portion of itemsets than  $nr(G; \mathcal{A})$  that, in turn, is anti-monotonic in more queries than  $nr(G; \mathcal{C})$ .

## 6 Conclusions

We have given a definition of a measure for ranking itemsets. The idea is to predict the frequency of an itemset from the frequencies of its sub-itemsets and measure the deviation between the actual frequency and the prediction. The more the itemset deviates from the prediction, the more it is significant. We estimated the frequencies using Maximum entropy and we used Kullback–Leibler divergence to measure the deviation. In the general case, the measure can be computed in  $O(2^{|G|})$  time, where  $|G|$  is the size of the itemset needed to be ranked, however, the measures  $r(G; \mathcal{T}^*)$  and  $r(G; \mathcal{I})$  can be computed in polynomial time.

We introduced two flexible rank measures  $r(G; \mathcal{T}^*)$  and  $r(G; \mathcal{F}^*)$ . The measure  $r(G; \mathcal{T}^*)$  can be solved by finding the optimal spanning tree in the mutual information matrix. For solving  $r(G; \mathcal{F}^*)$  we proposed a simple greedy approach.

A clear advantage of our approach to the previous methods is that the previous solutions calculate the deviation from the independence model whereas we are able to use the information available from the itemsets of larger size, and thus use more flexible models.

Our empirical results for real data show that the independence is too strict assumption: Almost all itemsets were significant according to  $r(G; \mathcal{I})$ . The results changed when we applied the more flexible models,  $r(G; \mathcal{C})$  and  $r(G; \mathcal{A})$ . We also observed an interesting type of overfitting: In some cases we obtain a better prediction if we do not use all the available information.

We showed that there is a little correlation between our measures and the other approaches. For instance, infrequent itemset may be significant and frequent itemset may be insignificant. We also observed that  $r(G; \mathcal{I})$  is monotonic for a large portion of itemsets, whereas  $r(G; \mathcal{C})$  and  $r(G; \mathcal{A})$  are anti-monotonic for a significant portion of itemsets.

**Acknowledgments** The author wishes to thank Gemma Garriga, Heikki Mannila, and Robert Gwadera for their comments.

### Appendix A. Asymptotic behaviour of the divergence

By asymptotic behaviour we mean the following: We assume that we have an ensemble of datasets  $D_i$  such that  $|D_i| \rightarrow \infty$ . We assume that  $G$  is non-derivable in each  $D_i$  and that the frequencies of  $\mathcal{F}_G$  are all equal.

Define  $N = |D|$  and  $M = |\mathcal{H}|$ . Let  $\mathbb{P}$  be the set of distributions satisfying the itemsets  $\mathcal{F}_G$ . It is easy to see that we can parameterize  $\mathbb{P}$  with frequencies of  $\mathcal{H}$ . In other words, let  $\mathcal{H} = \{H_1, \dots, H_M\}$ . Then for each  $p \in \mathbb{P}$ , there is a unique frequency vector  $\theta \in \mathbb{R}^M$  such that  $\theta_i = p(H_i = 1)$ . Let  $\Theta$  be the set of all possible frequency vectors. The set  $\Theta$  is a closed polytope—the vectors located on the boundary of  $\Theta$  corresponds to the distributions in which at least one entry is 0.

Let  $\theta^*$  be a frequency vector corresponding to the Maximum Entropy distribution  $p^*$ . We need to show that  $\theta^*$  is not a boundary vector. Assume the converse, then  $p^*$  must have  $p^*(\omega) = 0$  for some  $\omega$ . We know that this implies that  $p(\omega) = 0$  for all  $p \in \mathbb{P}$  [11, Theorem 3.1]. Let  $Y$  be the itemset containing the elements for which  $\omega$  has positive entries. This in turns (see [8]) implies that for each  $p \in \mathbb{P}$

$$p(G = 1) = \sum_{Y \subseteq Z \subseteq G} (-1)^{|G|-|Z|} p(Z = 1),$$

making  $G$  derivable and contradicting the statement.

Since  $\theta^*$  is an inner point of  $\Theta$ , let  $B \subset \Theta$  be an open ball around  $\theta^*$ . Assume that  $\theta \in B$ . By taking the expectation of the second-degree Taylor expansion of  $\log \frac{p(\omega; \theta^*)}{p(\omega; \theta)}$  around  $\theta$  we arrive to

$$-\text{KL}(\theta \parallel \theta^*) = \frac{1}{2} \Delta \theta^T \text{E}_\theta [H(\omega; \eta)] \Delta \theta,$$

where  $\Delta \theta = \theta^* - \theta$  and  $\eta$  is a vector lying between  $\theta$  and  $\theta^*$ , and  $H$  is the Hessian matrix of  $\log p(\omega; \eta)$ .

Let  $\theta_N$  be the frequencies of  $\mathcal{H}$  obtained from a dataset containing  $N$  points. According to 0-hypothesis we have  $\theta_N \rightsquigarrow \theta^*$  and  $\sqrt{N}(\theta_N - \theta^*) \rightsquigarrow N(0, \Sigma)$ , where  $\Sigma$  is a covariance matrix,

$$\Sigma_{ij} = p^*(H_i = 1, H_j = 1) - p^*(H_i = 1)p^*(H_j = 1).$$

If  $\theta_N \in B$ , we let  $\eta_N$  correspond to  $\eta$  in the Taylor expansion, otherwise we set  $\eta_N = 0$ . We can show that  $\eta_N \rightsquigarrow \theta^*$  [33, Theorem 2.7]. Consider a function

$$g(a, b, c, d) = \begin{cases} -a^T \text{E}_c [H(\omega; b)] a & c \in B \\ (2/d) \text{KL}(c \parallel \theta^*) & c \notin B \end{cases}.$$

This function is continuous in  $(\mathbb{R}^M, \theta^*, \theta^*, 0)$ . Hence, we can apply continuous map theory [33, Theorem 2.3] to obtain that

$$2N \text{KL}(\theta_N \parallel \theta^*) = g\left(\sqrt{N}(\theta_N - \theta^*), \eta_N, \theta_N, \frac{1}{N}\right) \rightsquigarrow -X^T \text{E}_{\theta^*} [H(\omega; \theta^*)] X,$$

where  $X$  is a random variable distributed as  $N(0, \Sigma)$ . We know that  $\text{E}_{\theta^*} [H(\omega; \theta^*)] = -\Sigma^{-1}$  [23, Lemma 4.11]. Theorem follows since  $X^T \Sigma^{-1} X$  is distributed as  $\chi^2$  with  $M$  degrees of freedom [33, Lemma 17.1].

## References

1. Aggarwal CC, Yu PS (1998) A new framework for itemset generation. In: PODS '98: Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems. ACM Press, New York, pp 18–24
2. Agrawal R, Imielinski T, Swami AN (1993) Mining association rules between sets of items in large databases. In: Buneman P, Jajodia S (eds) Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data. Washington, D.C., pp 207–216
3. Agrawal R, Mannila H, Srikant R, Toivonen H, Verkamo AI (1996) Fast discovery of association rules. In: Fayyad U, Piatetsky-Shapiro G, Smyth P, Uthurusamy R (eds) Advances in Knowledge Discovery and Data Mining. AAAI Press/The MIT Press, Cambridge pp 307–328
4. Boulicaut J-F, Bykowski A, Rigotti C (2000) Approximation of frequency queries by means of free-sets. In: Principles of Data Mining and Knowledge Discovery, pp 75–85
5. Brijs T, Swinnen G, Vanhoof K, Wets G (1999) Using association rules for product assortment decisions: a case study. In: Knowledge Discovery and Data Mining. ACM, New York, pp 254–260
6. Brin S, Motwani R, Silverstein C (1997) Beyond market baskets: Generalizing association rules to correlations. In: Peckham J (ed) SIGMOD 1997, Proceedings ACM SIGMOD International Conference on Management of Data. ACM Press, New York pp 265–276
7. Brin S, Motwani R, Ullman JD, Tsur S (1997) Dynamic itemset counting and implication rules for market basket data. In: SIGMOD 1997, Proceedings ACM SIGMOD International Conference on Management of Data. pp 255–264
8. Calders T, Goethals B (2002) Mining all non-derivable frequent itemsets. In: Proceedings of the 6th European Conference on Principles and Practice of Knowledge Discovery in Databases
9. Chow C, Liu C (1968) Approximating discrete probability distributions with dependence trees. IEEE Trans Inf Theory 14(3):462–467
10. Cooper G (1990) The computational complexity of probabilistic inference using bayesian belief networks. Artif Intell 42(2–3):393–405
11. Csiszár I (1975) I-divergence geometry of probability distributions and minimization problems. Ann Prob 3(1):146–158
12. Darroch J, Ratchli D (1972) Generalized iterative scaling for log-linear models. Ann Math Stat 43(5):1470–1480
13. Dong G, Li J (1999) Efficient mining of emerging patterns: Discovering trends and differences. In: Knowledge Discovery and Data Mining, pp 43–52
14. DuMouchel W, Pregibon D (2001) Empirical bayes screening for multi-item associations. In: Knowledge Discovery and Data Mining, pp 67–76
15. Fortelius M (2005) Neogene of the old world database of fossil mammals (NOW). University of Helsinki, <http://www.helsinki.fi/science/now/>
16. Fortelius M, Gionis A, Jernvall J, Mannila H (2006) Spectral ordering and biochronology of european fossil mammals paleobiology. Paleobiology 32(2):206–214
17. Gallo A, Bie TD, Christianini N (2007) Mini: Mining informative non-redundant itemsets. In: 11th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD), pp 438–445
18. Geurts K, Wets G, Brijs T, Vanhoof K (2003) Profiling high frequency accident locations using association rules. In: Proceedings of the 82nd Annual Transportation Research Board, Washington DC. (USA), January 12–16
19. Heikinheimo H, Hinkkanen E, Mannila H, Mielikäinen T, Seppänen JK (2007) Finding low-entropy sets and trees from binary data. In: Knowledge Discovery and Data Mining
20. Jaroszewicz S, Simovici DA (2002) Pruning redundant association rules using maximum entropy principle. In: Advances in Knowledge Discovery and Data Mining, 6th Pacific-Asia Conference, PAKDD'02, pp 135–147
21. Jiroušek R, Přeštil S (1995) On the effective implementation of the iterative proportional fitting procedure. Comput Stat Data Anal 19:177–189
22. Kohavi R, Brodley C, Frasca B, Mason L, Zheng Z (2000) KDD-Cup 2000 organizers report: Peeling the onion. SIGKDD Explorat 2(2):86–98
23. Kullback S (1968) Information Theory and Statistics. Dover Publications, Inc., New York
24. Mannila H, Mielikäinen T (2003) The pattern ordering problem. In: Principles of Data Mining and Knowledge Discovery, pp 327–338
25. Norén GN, Bate A, Edwards IR (2007) Extending the methods used to screen the who drug safety database towards analysis of complex associations and improved accuracy for rare events. Stat Med 25:3740–3757
26. Omiecinski ER (2003) Alternative interest measures for mining associations in databases. IEEE Trans Knowl Data Eng 15(1):57–69

27. Pasquier N, Bastide Y, Taouil R, Lakhal L (1999) Discovering frequent closed itemsets for association rules. *Lecture Notes in Computer Science* 1540:398–416
28. Pavlov D, Mannila H, Smyth P (2003) Beyond independence: Probabilistic models for query approximation on binary transaction data. *IEEE Trans Knowl Data Eng* 15(6):1409–1421
29. Piatetsky-Shapiro G (1991) Discovery, analysis, and presentation of strong rules. In: *Knowledge Discovery in Databases*. AAAI/MIT Press, New York, pp 229–248
30. Tatti N (2006a) Computational complexity of queries based on itemsets. *Inf Process Lett* pp 183–187
31. Tatti N (2006) Safe projections of binary data sets. *Acta Inf* 42(8–9):617–638
32. Tatti N (2007) Maximum entropy based significance of itemsets. In: *Proceedings of Seventh IEEE International Conference on Data Mining (ICDM 2007)*, pp 312–321
33. van der Vaart AW (1998) *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge
34. Webb GI (2006) Discovering significant rules. In: *Knowledge discovery and data mining*, pp 434–443

## Author Biography



**Nikolaj Tatti** is currently a Researcher at Basic Research Unit of Helsinki Institute for Information Technology, and also a Graduate Student at Department of Information and Computer Science of Helsinki University of Technology. His research interests include algorithms, statistics, and data mining.