

Technical Report

Department of Computer Science
and Engineering
University of Minnesota
4-192 EECS Building
200 Union Street SE
Minneapolis, MN 55455-0159 USA

TR 04-032

HARMONY: Efficiently Mining the Best Rules for Classification

Jianyong Wang and George Karypis

September 16, 2004

Report Documentation Page

Form Approved
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

| | | | | | |
|---|------------------------------------|-------------------------------------|----------------------------|--|---------------------------------|
| 1. REPORT DATE 16 SEP 2004 | | 2. REPORT TYPE | | 3. DATES COVERED - | |
| 4. TITLE AND SUBTITLE HARMONY: Efficiently Mining the Best Rules for Classification | | | | 5a. CONTRACT NUMBER | |
| | | | | 5b. GRANT NUMBER | |
| | | | | 5c. PROGRAM ELEMENT NUMBER | |
| 6. AUTHOR(S) | | | | 5d. PROJECT NUMBER | |
| | | | | 5e. TASK NUMBER | |
| | | | | 5f. WORK UNIT NUMBER | |
| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Army Research Laboratory, 2800 Powder Mill Road, Adelphi, MD, 20783-1197 | | | | 8. PERFORMING ORGANIZATION REPORT NUMBER | |
| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | | | | 10. SPONSOR/MONITOR'S ACRONYM(S) | |
| | | | | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) | |
| 12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited | | | | | |
| 13. SUPPLEMENTARY NOTES The original document contains color images. | | | | | |
| 14. ABSTRACT see report | | | | | |
| 15. SUBJECT TERMS | | | | | |
| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES 18 | 19a. NAME OF RESPONSIBLE PERSON |
| a. REPORT unclassified | b. ABSTRACT unclassified | c. THIS PAGE unclassified | | | |

HARMONY: Efficiently Mining the Best Rules for Classification *

Jianyong Wang and George Karypis

Department of Computer Science, Digital Technology Center, & Army HPC Research Center
University of Minnesota, Minneapolis, MN 55455
{jianyong, karypis}@cs.umn.edu

Abstract

Many studies have shown that rule-based classification algorithms perform well in classifying categorical and sparse high-dimensional databases. However, a fundamental limitation with many rule-based classifiers is that they find the classification rules in a coarse-grained manner. They usually use heuristic methods to prune the search space, and select the rules based on the sequential database covering paradigm. Thus, the so-mined rules may not be the globally best rules for some instances in the training database. To make worse, these algorithms fail to fully exploit some more effective search space pruning methods in order to scale to large databases.

In this paper we propose a new classifier, HARMONY, which directly mines the final set of classification rules. HARMONY uses an instance-centric rule-generation approach in the sense that it can assure for each training instance, one of the highest-confidence rules covering this instance is included in the result set, which helps a lot in achieving high classification accuracy. By introducing several novel search strategies and pruning methods into the traditional frequent item-set mining framework, HARMONY also has high efficiency and good scalability. Our thorough performance study with some large text and categorical databases has shown that HARMONY outperforms many well-known classifiers in terms of both accuracy and efficiency, and scales well w.r.t. the database size.

*This work was supported in part by NSF CCR-9972519, EIA-9986042, ACI-9982274, ACI-0133464, and ACI-0312828; the Digital Technology Center at UMN; and by the Army HPC Research Center under the auspices of the Department of the Army, Army Research Laboratory (ARL) under Cooperative Agreement number DAAD19-01-2-0014. The content of which does not necessarily reflect the position or the policy of the government, and no official endorsement should be inferred. Access to research and computing facilities was provided by the Minnesota Supercomputing Institute.

1 Introduction

As one of the most fundamental data mining tasks, classification has been extensively studied and various types of classification algorithms have been proposed. Among which, one category is the rule-based classifiers [?, ?, ?, ?]. They build a model from the training database as a set of high-quality rules, which can be used to predict the class labels of unlabeled instances. Many studies have shown that rule-based classification algorithms perform very well in classifying both categorical databases [?, ?, ?, ?] and databases represented via sparse high-dimensional such as those arising in the context of document classification [?, ?].

Some traditional rule-based algorithms like FOIL [?], RIPPER [?], and CPAR [?] discover a set of classification rules one-rule-at-a-time and employ a sequential covering methodology to eliminate from the training set the positive instances that are covered by each newly discovered rule. This *rule induction* process is done in a greedy fashion as it employs various heuristics (e.g., information gain) to determine how each rule would be extended. Due to this heuristic rule-induction process and the sequential covering framework, the final set of discovered rules are not guaranteed to be the best possible. For example, due to the removal of some training instances, the information gain is computed based on the incomplete information; thus, the variable (or literal) chosen by these algorithms to extend the current rule will be no longer the globally optimal one. Moreover, for multi-class problems, these algorithms need to be applied multiple times, each time mining the rules for one class. If the training database is large and contains many classes, the algorithms will be inefficient.

Since the introduction of association rule mining [?], many association-based (or related) classifiers have been proposed [?, ?, ?, ?, ?, ?, ?, ?, ?]. Some typical examples like CBA [?] and CMAR [?] adopt efficient association rule mining algorithms (e.g., Apri-

ori [?] and FP-growth [?]) to first mine a large number of high-confidence rules satisfying a user-specified minimum support and confidence thresholds and then use various sequential-covering-based schemes to select from them a set of high-quality rules to be used for classification. Since these schemes defer the selection step only after a large intermediate set of high-confidence rules have been identified, they tend to achieve somewhat better accuracy than the traditional heuristic rule induction schemes [?]. However, the drawback of these approaches is that the number of initial rules is usually extremely large, significantly increasing the rule discovery and selection time.

In this paper we propose a new classification algorithm, HARMONY¹, which can overcome the problems of both the rule-induction-based and the association-rule-based algorithms. HARMONY directly mines for each training instance one of the highest confidence classification rules that it supports and satisfies a user-specified minimum support constraint, and builds the classification model from the union of these rules over the entire set of instances. Thus HARMONY employs an *instance-centric* rule generation framework and is guaranteed to find and include the best possible rule for each training instance. Moreover, since each training instance usually supports many of the discovered rules, the overall classifier can better generalize to new instances and thus achieve better classification performance.

To achieve high computational efficiency, HARMONY mines the classification rules for all the classes simultaneously and directly mines the final set of classification rules by pushing deeply some effective pruning methods into the projection-based frequent itemset mining framework. All these pruning methods preserve the completeness of the resulting rule-set in the sense that they only remove from consideration rules that are guaranteed not to be of high quality. We have performed numerous performance studies with various databases and shown that HARMONY can achieve better accuracy while maintaining high efficiency.

The rest of the paper is organized as follows. Section ?? introduces some basic definitions and notations. Section ?? describes the problem formulation. Section ?? introduces some related work. Section ?? discusses in detail the HARMONY algorithm and some extensions to the algorithm. The thorough performance study is presented in Section ?. Finally, the paper concludes with Section ?.

¹ HARMONY stands for **H**igh confidence **A**ssociation **R**ule Mining **f**or **i**nstance-centric **c**lassif**Y**ing.

2 Notations and Definitions

A *training database* $TrDB$ is a set of training instances², where each training instance, denoted as a triple $\langle tid, X, cid \rangle$, contains a set of items (i.e., X) and is associated with a unique training instance identifier tid , and a class identifier $cid \in \{c_1, c_2, \dots, c_k\}$ (A class identifier is also called a class label, and we assume there are totally k distinct class labels in $TrDB$). Table ?? illustrates an example training database, which contains totally eight instances and two classes. Let $I = \{i_1, i_2, \dots, i_n\}$ be the complete set of distinct items appearing in $TrDB$. An *itemset* Y is a non-empty subset of I and is called an *l-itemset* if it contains l items. An itemset $\{x_1, \dots, x_l\}$ is also denoted by $x_1 \cdots x_l$. A training instance $\langle tid, X, cid \rangle$ is said to *contain* itemset Y if $Y \subseteq X$. The number of instances in $TrDB$ containing itemset Y is called the (absolute) *support* of itemset Y , denoted by sup_Y . The number of instances containing itemset Y and associated with a class label c_i (where $i \in \{1, 2, \dots, k\}$) is called the support of $Y \cup \{c_i\}$, denoted by $sup_Y^{c_i}$. A classification rule has the form: ' $Y \rightarrow c_i : sup_Y^{c_i}, conf_Y^{c_i}$ ', where Y is called the body, c_i the head, $sup_Y^{c_i}$ the support, and $conf_Y^{c_i} = \frac{sup_Y^{c_i}}{sup_Y}$ the confidence of the rule, respectively. In addition, we use $|TrDB|$ to denote the number of instances in database $TrDB$, and for brevity, we sometimes use the instance identifier tid to denote an instance $\langle tid, X, cid \rangle$.

Table 1. An example training database $TrDB$.

| Instance identifier | Set of items | Class identifier |
|---------------------|--------------|------------------|
| 01 | a, c, e, g | 1 |
| 02 | b, d, e, f | 0 |
| 03 | d, e, f | 0 |
| 04 | a, b, c, e | 1 |
| 05 | a, c, e | 1 |
| 06 | b, d, e | 0 |
| 07 | a, b, e | 1 |
| 08 | a, b, d, e | 0 |

Given a minimum support threshold, min_sup , an itemset Y is *frequent* if $sup_Y \geq min_sup$. A frequent itemset Y supported by any training instance $\langle t_j, X_j, c_i \rangle$ ($1 \leq j \leq |TrDB|$ and $1 \leq i \leq k$) is also called a frequent covering itemset of instance t_j , and ' $Y \rightarrow c_i : sup_Y^{c_i}, conf_Y^{c_i}$ ' is called a frequent covering

²Note there may exist a test database, which is in the same form as the training database and is used to evaluate the performance of a classifier. We denote it by $TeDB$.

rule of instance t_j ³. Among the frequent covering rules of any instance t_j , those with the highest confidence are called the *Highest Confidence Covering Rules* w.r.t. instance t_j . We denote a Highest Confidence Covering Rule w.r.t. instance t_j by $HCCR_{t_j}$, and use $HCCR_{t_j}^{sup}$ and $HCCR_{t_j}^{conf}$ to denote its support and confidence.

3 Problem Definition

The goal of this paper is to design an accurate and efficient rule-based classifier with good scalability, which should be able to overcome the problems of both the traditional rule-based and the recently proposed association-based classifiers. As mentioned in Section ??, instead of using the sequential database covering to select the rules, our solution mines a set of high quality rules in an instance-centric manner and can assure that at least one of the highest confidence frequent covering rules (if there is any) w.r.t. any training instance is included in the final result set of classification rules.

Specifically, given a training database $TrDB$ and a minimum support threshold min_sup , the problem of this study is to find one of the highest confidence frequent covering rules for each of the training instances in $TrDB$, and build a classifier from these classification rules. Note the input training database must be in the form that is consistent with the corresponding definition in Section ??, otherwise, the training database should be first converted to that form. For example, a numerical database should be first discretized into a categorical one in order to use HARMONY to build the model. In addition, although this study mainly focuses on mining any one of the highest confidence frequent covering rules for each training instance, it is straightforward to revise HARMONY to mine the complete set of highest confidence frequent covering rules or K highest confidence frequent covering rules for each training instance.

4 Related Work

There are two classes of algorithms that are directly related to this work. One is the traditional rule-induction-based methods and the other is the recently proposed association-rule-based methods. Both of these classes share the same idea of trying to find

³Note in this paper by a ‘frequent rule’ we mean its rule body is frequent. This is different from the traditional definition of an association rule, which requires the support of the entire rule is no lower than a minimum support. However, it is straightforward to adapt HARMONY to the traditional definition as discussed in Section ??.

a set of classification rules to build the model. The rule-induction-based classifiers like C4.5 [?], FOIL [?], RIPPER [?], and CPAR [?] use various heuristics such as information gain (including Foil gain) and gini index to identify the best variable (or literal) by which to grow the current rule, and many of them follow a sequential database covering paradigm to speed up rule induction. The association-based classifiers adopt another approach to find the set of classification rules. They first use some efficient association rule mining algorithms to discover the complete (or a large intermediate) set of association rules, from which the final set of classification rules can be chosen based on various types of sequential database covering techniques. Some typical examples of association-based methods include CBA [?], CMAR [?], and ARC-BC [?].

In contrast to the rule-induction-based algorithms, HARMONY does not apply any heuristic pruning methods and the sequential database covering approach. Instead, it follows an instance-centric framework and mines the covering rules with the highest confidence for each instance, which can achieve better accuracy. At the same time, by maintaining the currently best rules for each training instance and pushing deeply several effective pruning methods into the projection-based frequent itemset mining framework [?, ?, ?], HARMONY directly mines the final set of classification rules, which avoids the time consuming rule generation and selection process used in several association-based classifiers [?, ?, ?].

The idea of directly mining a set of high confidence classification rules is similar to those in [?, ?]. The author of [?] investigated a brute-force technique for mining the set of high-confidence classification rules, and proposed several effective pruning strategies to control the combinatorial explosion in the number of rule candidates. The FARMER algorithm [?] finds the interesting rule groups for microarray databases. It mines the rules in a row enumeration space, and fully exploits some effective pruning methods to prune the search space based on the user-specified constraints like minimum support, confidence, and chi-square. Unlike [?, ?], HARMONY does not need the user to specify the minimum confidence and/or chi-square. Instead, it mines for each training instance one of the highest confidence frequent rules that it covers. In addition, by maintaining the currently best classification rules for each instance, HARMONY is able to incorporate some new pruning methods under the unpromising item (or conditional database) pruning framework, which has been proven very effective in pushing deeply the length-decreasing support constraint or tough block constraints into closed itemset mining [?, ?].

Contributions. We summarize the contributions of this paper as follows.

1. We proposed an instance-centric paradigm in mining the highest confidence covering rules for each training instance, which can be used to build an accurate classification model.
2. We explored new search strategies and pruning methods to effectively reduce the search space.
3. A new classifier, HARMONY, was designed, and a thorough performance study with various large text and categorical databases has shown that HARMONY is very accurate and efficient compared to many well-known classifiers.

5 HARMONY: An Instance-Centric Classifier

In this section, we will describe in detail the HARMONY algorithm. We first elaborate on how to adapt the traditional projection-based frequent itemset mining framework to efficiently enumerate the classification rules, then we focus on how to push deeply some effective pruning methods into the rule enumeration framework. Finally we will give the whole algorithm.

5.1 Classification Rule Enumeration

The projection-based itemset enumeration framework has been widely used in many frequent itemset mining algorithms [?, ?, ?], and will be used by HARMONY as the basis in enumerating the classification rules. Given a training database $TrDB$ and a minimum support min_sup , HARMONY first computes the frequent items by scanning $TrDB$ once, and sorts them to get a list of frequent items (denoted by f_list) according to a certain ordering scheme. Assume the min_sup is 3 and the lexicographical ordering is the default ordering scheme, the f_list computed from Table ?? is $\{a, b, c, d, e\}$. HARMONY applies the *divide-and-conquer* method plus the *depth-first search strategy*. In our example, HARMONY first mines the rules whose body contains item ‘a’, then mines the rules whose body contains ‘b’ but no ‘a’, ..., and finally mines the rules whose body contains only ‘e’. In mining the rules with item ‘a’, item ‘a’ is treated as the current prefix, and its conditional database (denoted by $TrDB|_a$) is built and the *divide-and-conquer* method is applied recursively with the depth-first search strategy. To build conditional database $TrDB|_a$, HARMONY first identifies the instances in $TrDB$ containing ‘a’ and removes the infrequent items, then sorts the left items in each instance

according to the f_list order, finally $TrDB|_a$ is built as $\{\langle 01, ce, 1 \rangle, \langle 04, bce, 1 \rangle, \langle 05, ce, 1 \rangle, \langle 07, be, 1 \rangle, \langle 08, be, 0 \rangle\}$ (infrequent items ‘d’ and ‘g’ are removed). Following the *divide-and-conquer* method, HARMONY first mines the rules with prefix ‘ab’, then mines rules with prefix ‘ac’ but no ‘b’, and finally mines rules with prefix ‘ae’ but no ‘b’ nor ‘c’.

During the mining process, when HARMONY gets a new prefix, it will generate a set of classification rules w.r.t. the training instances covered by the prefix. For each training instance, it always maintains one of its currently highest confidence rules mined so far. Assume the current prefix P is ‘a’ (i.e., $P='a'$). As shown in the above example, P covers five instances with $tids$ 01, 04, 05, 07, and 08. HARMONY computes the covering rules according to the class distribution w.r.t. the prefix P . In this example, $sup_P=5$, $sup_P^0=1$, $sup_P^1=4$, and HARMONY generates two classification rules:

$$\begin{aligned} \text{Rule 1: } & a \rightarrow 0 : 1, \frac{1}{5} \\ \text{Rule 2: } & a \rightarrow 1 : 4, \frac{4}{5} \end{aligned}$$

Rule 1 covers the instance with tid 08, while Rule 2 covers the instances with $tids$ 01, 04, 05 and 07. Up to this point, we have $HCCR_{01} = HCCR_{04} = HCCR_{05} = HCCR_{07} = \text{Rule 2}$, and $HCCR_{08} = \text{Rule 1}$.

5.1.1 Ordering of the Local Items

In the above rule enumeration process, we used the lexicographical ordering as an illustration to sort the set of local frequent items in order to get the f_list . Many frequent itemset mining algorithms either adopt item support descending order [?] or support ascending order [?] as the ordering scheme. However, because we are interested in the highest confidence rules w.r.t. the training instances, both the support descending order and ascending order may not be the most efficient and effective ways. As a result, we propose the following three new ordering schemes as the alternatives.

Let the current prefix be P , its support be sup_P , the support and confidence of the classification rule w.r.t. prefix P and class label c_i , ‘ $P \rightarrow c_i$ ’, be $sup_P^{c_i}$ and $conf_P^{c_i}$, respectively, the set of local frequent items be $\{x_1, x_2, \dots, x_m\}$, the number of prefix P ’s conditional instances containing item x_j ($1 \leq j \leq m$) and associated with class label c_i ($1 \leq i \leq k$) be $sup_{P \cup \{x_j\}}^{c_i}$, and the support of $P \cup \{x_j\}$ be $sup_{P \cup \{x_j\}} = \sum_{i=1}^k sup_{P \cup \{x_j\}}^{c_i}$.

Maximum confidence descending order. Given a local item x_j ($1 \leq j \leq m$) w.r.t. P , we can compute k rules with body $P \cup \{x_j\}$, among which, the i -th rule with rule head c_i is:

$$P \cup \{x_j\} \rightarrow c_i : sup_{P \cup \{x_j\}}^{c_i}, \frac{sup_{P \cup \{x_j\}}^{c_i}}{sup_{P \cup \{x_j\}}}$$

The highest confidence among the k rules with body $P \cup \{x_j\}$ is called the maximum confidence of local item x_j , and is defined as the following:

$$\frac{\max_{\forall i, 1 \leq i \leq k} \text{sup}_{P \cup \{x_j\}}^{c_i}}{\text{sup}_{P \cup \{x_j\}}} \quad (1)$$

To mine the highest confidence covering rules as quickly as possible, a good heuristic is to sort the local frequent items in their maximum confidence descending order.

Entropy ascending order. The widely used entropy to some extent measures the purity of a cluster of instances. If the entropy of the set of instances containing $P \cup \{x_j\}$ ($1 \leq j \leq m$) is small, it is highly possible to generate some high confidence rules with body $P \cup \{x_j\}$. Thus another good ordering heuristic is to rank the set of local frequent items in their entropy ascending order, and the entropy w.r.t. item x_j is defined as follows [?]:

$$-\frac{1}{\log k} \sum_{i=1}^k \left(\frac{\text{sup}_{P \cup \{x_j\}}^{c_i}}{\text{sup}_{P \cup \{x_j\}}} \right) \log \left(\frac{\text{sup}_{P \cup \{x_j\}}^{c_i}}{\text{sup}_{P \cup \{x_j\}}} \right) \quad (2)$$

Correlation coefficient ascending order. Both the maximum confidence descending order and entropy ascending order do not consider the class distribution of the conditional database w.r.t. prefix P , which may cause some problems in some cases. Let us see an example. Assume the number of class labels $k=2$, $\text{sup}_P^{c_1} = 12$, and $\text{sup}_P^{c_2} = 6$, then we can get two rules with body P as follows:

$$\text{Rule 3: } P \rightarrow c_1 : 12, \frac{12}{18}$$

$$\text{Rule 4: } P \rightarrow c_2 : 6, \frac{6}{18}$$

Suppose there are two local items, x_1 and x_2 , and $\text{sup}_{P \cup \{x_1\}}^{c_1} = 2$, $\text{sup}_{P \cup \{x_1\}}^{c_2} = 1$, $\text{sup}_{P \cup \{x_2\}}^{c_1} = 1$, and $\text{sup}_{P \cup \{x_2\}}^{c_2} = 2$. According to Equation ?? and Equation ??, the maximum confidence and entropy w.r.t. item x_1 are equal to the corresponding maximum confidence and entropy w.r.t. x_2 . Thus we cannot determine which one of x_1 and x_2 should be ranked higher. However, because the conditional database $\text{TrDB}|_{P \cup \{x_1\}}$ has the same class distribution as conditional database $\text{TrDB}|_P$, we cannot generate rules with body $P \cup \{x_1\}$ and a confidence higher than those with body P (i.e., Rule 3 and Rule 4). The two rules with body $P \cup \{x_1\}$ are shown as the following.

$$\text{Rule 5: } P \cup \{x_1\} \rightarrow c_1 : 2, \frac{2}{3}$$

$$\text{Rule 6: } P \cup \{x_1\} \rightarrow c_2 : 1, \frac{1}{3}$$

If we examine the rules generated from prefix item-set $P \cup \{x_2\}$ as shown in Rule 7 and Rule 8, we can see Rule 8 has higher confidence than Rule 4, and can be used to replace Rule 4 for the instances covered by

Rule 8. In this case, item x_2 should be ranked before item x_1 .

$$\text{Rule 7: } P \cup \{x_2\} \rightarrow c_1 : 1, \frac{1}{3}$$

$$\text{Rule 8: } P \cup \{x_2\} \rightarrow c_2 : 2, \frac{2}{3}$$

This example suggests that the more similar the class distribution between conditional databases $\text{TrDB}|_P$ and $\text{TrDB}|_{P \cup \{x_j\}}$ ($1 \leq j \leq m$), the lower is the possibility to generate higher confidence rules from $\text{TrDB}|_{P \cup \{x_j\}}$. Because the correlation coefficient is a good metric in measuring the similarity between two vectors (the larger the coefficient, the more similar the two vectors), it can be used to rank the local items. In HARMONY, the correlation coefficient ascending order is adopted to sort the local items.

Let $\overline{\text{sup}}_P$ be $\frac{1}{k} \sum_{i=1}^k \text{sup}_P^{c_i}$, $\overline{\text{sup}}_{P \cup \{x_j\}}$ be $\frac{1}{k} \sum_{i=1}^k \text{sup}_{P \cup \{x_j\}}^{c_i}$, σ_P be $\sqrt{\frac{1}{k} \sum_{i=1}^k (\text{sup}_P^{c_i})^2 - \overline{\text{sup}}_P^2}$, $\sigma_{P \cup \{x_j\}}$ be $\sqrt{\frac{1}{k} \sum_{i=1}^k (\text{sup}_{P \cup \{x_j\}}^{c_i})^2 - \overline{\text{sup}}_{P \cup \{x_j\}}^2}$, the correlation coefficient between prefix P and $P \cup \{x_j\}$ ($1 \leq j \leq m$) is defined as follows.

$$\frac{\frac{1}{k} \sum_{i=1}^k (\text{sup}_P^{c_i} \times \text{sup}_{P \cup \{x_j\}}^{c_i} - \overline{\text{sup}}_P \times \overline{\text{sup}}_{P \cup \{x_j\}})}{\sigma_P \times \sigma_{P \cup \{x_j\}}} \quad (3)$$

5.2 Search Space Pruning

Unlike the association-based algorithms, HARMONY directly mines the final set of classification rules. By maintaining the current highest confidence among the covering rules for each training instance during the mining process, some effective pruning methods can be proposed to improve the algorithm efficiency.

5.2.1 Support Equivalence Item Elimination

Given the current prefix P , among its set of local frequent items $\{x_1, x_2, \dots, x_m\}$, some may have the same support as P . We call them support equivalence items and can be safely pruned according to the following Lemma ??.

Lemma 1 (*Support equivalence item pruning*)

Any local item x_j w.r.t. prefix P can be safely pruned if it satisfies $\text{sup}_{P \cup \{x_j\}} = \text{sup}_P$.

Proof. Because $\text{sup}_{P \cup \{x_j\}} = \text{sup}_P$ holds, $\text{TrDB}|_P$ and $\text{TrDB}|_{P \cup \{x_j\}}$ contain the same set of conditional instances; thus, their class distributions are also the same and the following equation must hold:

$$\forall i, 1 \leq i \leq k, \text{sup}_{P \cup \{x_j\}}^{c_i} = \text{sup}_P^{c_i}$$

Given any itemset, Y , which can be used to extend P (Y can be empty), can also be used to extend $P \cup \{x_j\}$, and the following must hold:

$$\forall i, 1 \leq i \leq k, \sup_{P \cup \{x_j\} \cup Y}^{c_i} = \sup_{P \cup Y}^{c_i}$$

We can further have the following equation:

$$\forall i, 1 \leq i \leq k, \frac{\sup_{P \cup \{x_j\} \cup Y}^{c_i}}{\sup_{P \cup \{x_j\} \cup Y}} = \frac{\sup_{P \cup Y}^{c_i}}{\sup_{P \cup Y}}$$

This means the confidence of the rule ' $P \cup \{x_j\} \cup Y \rightarrow c_i$ ' is equal to the confidence of the rule ' $P \cup Y \rightarrow c_i$ ', and we cannot generate higher confidence rules from prefix $P \cup \{x_j\} \cup Y$ in comparison with the rules with body $P \cup Y$. Thus item x_j can be safely pruned. \square

Note $P \cup Y$ is a subset of $P \cup \{x_j\} \cup Y$, by pruning item x_j , we prefer the more generic classification rules. A similar strategy was adopted in [?, ?].

5.2.2 Unpromising Item Elimination

Given the current prefix P , any one of its local frequent items, x_j ($1 \leq j \leq m$), any itemset Y that can be used to extend $P \cup \{x_j\}$ (where Y can be empty and $P \cup \{x_j\} \cup Y$ is frequent), and any class label c_i ($1 \leq i \leq k$), the following equation must hold:

$$\begin{aligned} \text{conf}_{P \cup \{x_j\} \cup Y}^{c_i} &= \frac{\sup_{P \cup \{x_j\} \cup Y}^{c_i}}{\sup_{P \cup \{x_j\} \cup Y}} \leq \frac{\sup_{P \cup \{x_j\} \cup Y}^{c_i}}{\text{min_sup}} \\ &\leq \frac{\sup_{P \cup \{x_j\}}^{c_i}}{\text{min_sup}} \end{aligned}$$

Because $\text{conf}_{P \cup \{x_j\} \cup Y}^{c_i} \leq 1$ also holds, we have the following equation:

$$\text{conf}_{P \cup \{x_j\} \cup Y}^{c_i} \leq \min\left\{1, \frac{\sup_{P \cup \{x_j\}}^{c_i}}{\text{min_sup}}\right\} \quad (4)$$

Lemma 2 (Unpromising item pruning) For any conditional instance $\langle t_l, X_l, c_i \rangle \in \text{TrDB}|_{P \cup \{x_j\}}$ ($\forall l, 1 \leq l \leq |\text{TrDB}|_{P \cup \{x_j\}}$, and $1 \leq i \leq k$), if the following always holds, item x_j is called an unpromising item and can be safely pruned.

$$\text{HCCR}_{t_l}^{\text{conf}} \geq \min\left\{1, \frac{\sup_{P \cup \{x_j\}}^{c_i}}{\text{min_sup}}\right\} \quad (5)$$

Proof. By combining Equation ?? and Equation ?? we get that for any itemset Y (Y can be empty) the following must hold:

$$\text{conf}_{P \cup \{x_j\} \cup Y}^{c_i} \leq \text{HCCR}_{t_l}^{\text{conf}}$$

This means that any rule mined by growing prefix $P \cup \{x_j\}$ will have a confidence that is no greater than the current highest confidence covering rules (with the same rule head) of any conditional instance in $\text{TrDB}|_{P \cup \{x_j\}}$; thus, item x_j can be safely pruned. \square

5.2.3 Unpromising Conditional Database Elimination

Given the current prefix P , any itemset Y (where Y can be empty and $P \cup Y$ is frequent), any class label c_i ($1 \leq i \leq k$), the confidence of rule ' $P \cup Y \rightarrow c_i$ ', $\text{conf}_{P \cup Y}^{c_i}$, must satisfy the following equation:

$$\text{conf}_{P \cup Y}^{c_i} = \frac{\sup_{P \cup Y}^{c_i}}{\sup_{P \cup Y}} \leq \frac{\sup_{P \cup Y}^{c_i}}{\text{min_sup}} \leq \frac{\sup_P^{c_i}}{\text{min_sup}}$$

In addition, because $\text{conf}_{P \cup Y}^{c_i} \leq 1$ also holds, we have the following equation:

$$\text{conf}_{P \cup Y}^{c_i} \leq \min\left\{1, \frac{\sup_P^{c_i}}{\text{min_sup}}\right\} \quad (6)$$

Lemma 3 (Unpromising conditional database pruning) For any conditional instance $\langle t_l, X_l, c_i \rangle \in \text{TrDB}|_P$ ($\forall l, 1 \leq l \leq |\text{TrDB}|_P$, and $1 \leq i \leq k$), if the following always holds, the conditional database $\text{TrDB}|_P$ can be safely pruned.

$$\text{HCCR}_{t_l}^{\text{conf}} \geq \min\left\{1, \frac{\sup_P^{c_i}}{\text{min_sup}}\right\} \quad (7)$$

Proof. By combining Equation ?? and Equation ?? we can get that for any itemset Y (Y can be empty) and $\forall l, 1 \leq l \leq |\text{TrDB}|_P$, $\langle t_l, X_l, c_i \rangle \in \text{TrDB}|_P$ ($1 \leq i \leq k$), the following must hold:

$$\text{conf}_{P \cup Y}^{c_i} \leq \text{HCCR}_{t_l}^{\text{conf}}$$

This means that any rule mined by growing prefix P will have a confidence that is no greater than the current highest confidence rules (with the same rule head) of any conditional instance in $\text{TrDB}|_P$; thus, the whole conditional database $\text{TrDB}|_P$ can be safely pruned. \square

ALGORITHM 1: HARMONY(TrDB , min_sup , t_i)

INPUT: (1) TrDB : a training database, (2) min_sup : a minimum support threshold, and (3) t_i : a new test instance.

OUTPUT: (1) HCCR : the set of the highest confidence frequent covering rules w.r.t. each instance in TrDB , (2) CM : a classification model, (3) PCL : the predicted class label(s) w.r.t. test instance t_i .

01. $\text{HCCR} \leftarrow \text{RULEMINER}(\text{TrDB}, \text{min_sup})$;
 02. $\text{CM} \leftarrow \text{BUILDMODEL}(\text{HCCR})$;
 03. $\text{PCL} \leftarrow \text{NEWINSTANCECLASSIFICATION}(\text{CM}, t_i)$.

5.3 The algorithm

After we described how to enumerate the classification rules, and how to design the local item ordering scheme and some effective search space pruning methods in order to accelerate the mining of the highest

confidence covering rules in terms of each training instance, we introduce the integrated HARMONY algorithm in this section.

The HARMONY algorithm is shown in ALGORITHM 1. It consists of three sub-algorithms: RULEMINER() takes as input the training database $TrDB$ and the minimum support min_sup , and outputs the set of highest confidence covering classification rules, $HCCR$; BUILDMODEL() takes $HCCR$ as input and outputs a classification model, CM ; NEWINSTANCECLASSIFICATION() classifies a new test instance ti using the model CM .

5.3.1 Classification Rule Generation

In Section ?? and Section ?? we introduced how to efficiently enumerate the classification rules under the *divide-and-conquer* and *depth-first search* paradigm, and proposed several pruning methods to speed up the enumeration of the highest confidence covering rules. By integrating the pruning methods with the rule enumeration, we get the classification rule generation algorithm, as shown in the RULEMINER() algorithm.

The RULEMINER() algorithm first initializes the highest confidence classification rules w.r.t. each training instance to empty (lines 01-02), then enumerates the classification rules by calling subroutine *ruleminer*($\emptyset, TrDB$) (line 03). Subroutine *ruleminer*() takes as input a prefix itemset pi and its corresponding conditional database cdb . For each conditional instance, it checks if a classification rule with higher confidence can be computed from the current prefix pi , if so, it replaces the corresponding instance's current highest confidence rule with the new rule (lines 04-07). It then finds the frequent local items by scanning cdb (line 08), prunes invalid items based on the *support equivalence item pruning* method and the *unpromising item pruning* method (lines 09-10). If the set of valid local items is empty or the whole conditional database cdb can be pruned based on the *unpromising conditional database pruning* method, it returns directly (lines 11-13). Otherwise, it sorts the left frequent local items according to the correlation coefficient ascending order (line 14), and grows the current prefix (line 16), builds the conditional database for the new prefix (line 17), and recursively calls itself to mine the highest confidence rules from the new prefix (line 18).

5.3.2 Building the Classification Model

After the set of highest confidence covering rules have been mined, it will be straightforward to build the classification model. HARMONY first groups the set of

highest confidence covering rules into k groups according to their rule heads (i.e., class labels), where k is the total number of distinct class labels in the training database. Within the same group of rules, HARMONY sorts the rules in their confidence descending order, and for the rules with the same confidence, sorts them in support descending order. In this way, HARMONY prefers the rules with higher confidence, and the rules with higher support if the confidence is the same. The BUILDMODEL algorithm is shown in ALGORITHM 1.2.

ALGORITHM 1.1: RULEMINER($TrDB, min_sup$)

INPUT: (1) $TrDB$: a training database, and (2) min_sup : a minimum support threshold.
 OUTPUT: (1) $HCCR$: the set of the highest confidence frequent covering rules w.r.t. each instance in $TrDB$.

01. for all $t_i \in TrDB$
 02. $HCCR_{t_i} \leftarrow \emptyset$;
 03. call **ruleminer**($\emptyset, TrDB$).
-

SUBROUTINE 1.1 : **ruleminer**(pi, cdb)

INPUT: (1) pi : a prefix itemset, and (2) cdb : the conditional database w.r.t. prefix pi .

04. if($pi \neq \emptyset$)
 05. for all $\langle t_i, X_i, c_j \rangle \in cdb$
 06. if($HCCR_{t_i}^{conf} < \frac{sup_{c_j}}{sup_{pi}}$)
 07. $HCCR_{t_i} \leftarrow \text{rule } 'pi \rightarrow c_j'$;
 08. $I \leftarrow \text{find_frequent_items}(cdb, min_sup)$;
 09. $S \leftarrow \text{support_equivalence_item_pruning}(I)$; $I \leftarrow I - S$;
 10. $S \leftarrow \text{unpromising_item_pruning}(I, cdb)$; $I \leftarrow I - S$;
 11. if($I \neq \emptyset$)
 12. if($\text{unpromising_conditional_database_pruning}(I, pi, cdb)$)
 13. return;
 14. $\text{correlation_coefficient_ascending_ordering}(I)$;
 15. for all $x \in I$ do
 16. $pi' \leftarrow pi \cup \{x\}$;
 17. $cdb' \leftarrow \text{build_cond_database}(pi', cdb)$;
 18. call **ruleminer**(pi', cdb');
-

ALGORITHM 1.2: BUILDMODEL($HCCR$)

INPUT: (1) $HCCR$: the set of highest confidence covering rules.
 OUTPUT: (1) CM : the classification model (i.e., k groups of ranked rules).

01. $\text{Cluster_rules_into_k_groups}(HCCR)$; //according to class label
 02. for each group of rules
 03. $\text{Sort_rules}()$; //in confidence and support descending order
-

ALGORITHM 1.3: NEWINSTANCECLASSIFICATION(CM, ti)

INPUT: (1) CM : the classification model, (2) ti : a test instance.
 OUTPUT: (1) PCL : a predicted class label (or a set of class labels).

01. for $j=1$ to k // CM_j : the j -th group of rules in CM
 // SCR_j : the score for ti computed from CM_j
 02. $SCR_j \leftarrow \text{ComputeScore}(CM_j, ti)$;
 03. $PCL \leftarrow \text{PredictClassLabel}(SCR)$.
-

5.3.3 New Instance Classification

After the classification model, CM , has been built, it can be used to classify a new test instance, ti , using ALGORITHM 1.3. HARMONY first computes a score w.r.t. ti for each group of rules in CM (lines 01-02), and predicts for ti a class label or a set of class labels if the underlying classification is a multi-class multi-label problem (i.e., each instance can be associated with several class labels). In HARMONY, the score for a certain group of rules is defined as the sum of the confidences of the covering rules w.r.t. ti (by a ‘covering rule’, we mean its rule body is a subset of ti). For a multi-class single-label classification problem, HARMONY simply chooses the class label with the highest score as the predicted class label. While for a multi-class multi-label problem, a *dominant factor*-based method [?] can be used to predict the class labels and works as follows. Given a user-specified *dominant factor* γ , let the class label with the highest score be c_{max} and the corresponding highest score w.r.t. test instance ti be $SCORE_{ti}^{c_{max}}$, then any class label whose corresponding score is no smaller than $SCORE_{ti}^{c_{max}} \times \gamma$ is a predicted class label for ti .

Although the *dominant factor*-based method works well in many cases, if the distribution of the class labels in the training database is imbalanced, the average confidence of each group of classification rules may be quite different from each other, this uniform *dominant factor*-based method will have some problems. A large *dominant factor* may lead to low recalls (i.e., the percentage of the total test instances for the given class label that are correctly classified) for the classes with low average rule confidences, while a small *dominant factor* can lead to low precisions (i.e., the percentage of predicted instances for the given class label that are correctly classified) for the classes with high average rule confidences. To overcome this problem, HARMONY adopts a *weighted dominant factor*-based method. Let the average confidence of the group of classification rules w.r.t. class label c_k be $conf_{c_k}^{avg}$, the score w.r.t. instance ti and class label c_k be $SCORE_{ti}^{c_k}$. Instance ti is predicted to belong to class c_k if it satisfies the following equation:

$$SCORE_{ti}^{c_k} \geq SCORE_{ti}^{c_{max}} \times \gamma \times \left(\frac{conf_{c_k}^{avg}}{conf_{c_{max}}^{avg}} \right)^\delta$$

Here, δ ($\delta \geq 0$) is called the *score differentia factor* and the larger the δ , the more the difference of the *weighted dominant factors* (i.e., $\gamma \times \left(\frac{conf_{c_k}^{avg}}{conf_{c_{max}}^{avg}} \right)^\delta$) among different class labels. It is set to 1 by default in HARMONY.

5.4 Extensions

5.4.1 Varying Support Threshold

The classification rule enumeration algorithm described in Section ?? assumes a uniform minimum support as an input, which may cause some problems for the unbalanced training databases. By an unbalanced training database, we mean the class distribution is not balanced, that is, some classes may contain a much larger number of instances than the other classes. If a large minimum support is used as input, the algorithm will encounter difficulties in mining high confidence rules for the small classes, while a small minimum support will lead to the overfitting problem for some large classes. This intuition suggests we should use different minimum supports for different size classes.

HARMONY provides two ways in specifying varying minimum supports for different classes. The first way allows the user to directly specify a minimum support for each class (in the following we will use min_sup_i to denote the minimum support of the i -th class). However, when there exist a lot of classes in the database, to specify a proper minimum support for each class is not an easy task. As a result, in the second option, HARMONY requires the user to provide a minimum support, min_sup , which corresponds to the minimum support of the smallest class, and it will automatically compute a minimum support for each class from min_sup and the class distribution. Let the number of training instances w.r.t. class c_i be $|c_i|$, then min_sup_i is computed as follows:

$$min_sup_i = min_sup \times \left(\frac{|c_i|}{\min_{j, 1 \leq j \leq k} |c_j|} \right)^\xi$$

Here, ξ ($\xi \geq 0$) is called the *support differentia factor*. In HARMONY, ξ is set to 0 by default, which can be used to compute a uniform minimum support for all the classes.

By using varying minimum support, Lemma ?? and Lemma ?? still applies, but we need to replace min_sup with min_sup_i in Equation ?? and Equation ?. For example, Equation ?? should be changed to the following form:

$$HCCR_{ti}^{conf} \geq \min \left\{ 1, \frac{sup_{P \cup \{x_j\}}^{c_i}}{min_sup_i} \right\}$$

In addition, to make the algorithm work, we also need to require the line 06 of SUBROUTINE 1.1 satisfy $sup_{pi} \geq min_sup_j$.

5.4.2 Mining K -Rules for Each Instance

A training instance may support multiple highest-confidence classification rules, but the above classification rule enumeration algorithm described in Section ?? only reports the first discovered one. Usually this arrangement can assure the set of final rules is large enough to build an accurate classifier in the case that the training database contains a large number of training instances. However, the set of final rules generated in this way may not be sufficient if the database is small. To overcome this problem, HARMONY provides an option to mine K highest-confidence rules w.r.t. a training instance if it supports multiple highest-confidence rules, where K is a user-specified parameter.

In order to mine K -rules for each instance, Equation ?? needs to be changed to the following form:

$$(HCCR_{t_l}^{conf} > \min\{1, \frac{sup_{P \cup \{x_j\}}^{c_i}}{min_sup}\}) \vee$$

$$(HCCR_{t_l}^{conf} = \min\{1, \frac{sup_{P \cup \{x_j\}}^{c_i}}{min_sup}\}) \wedge (n \geq K)$$

Here, n is the number of highest confidence classification rules discovered so far w.r.t. t_l .

Similarly, Equation ?? should have the following form:

$$(HCCR_{t_l}^{conf} > \min\{1, \frac{sup_P^{c_i}}{min_sup}\}) \vee$$

$$(HCCR_{t_l}^{conf} = \min\{1, \frac{sup_P^{c_i}}{min_sup}\}) \wedge (n \geq K)$$

In addition, it is evident that the condition of line 06 in SUBROUTINE 1.1 should be rewritten to the following form:

$$(HCCR_{t_l}^{conf} < \frac{sup_{p_i}^{c_j}}{sup_{p_i}}) \vee (HCCR_{t_l}^{conf} = \frac{sup_{p_i}^{c_j}}{sup_{p_i}}) \wedge (n < K)$$

5.4.3 Traditional Definition of a Frequent Rule

The above classification rule enumeration algorithm described in Section ?? can also be adapted to accord with the more traditional definition of an association rule, that is, instead of only requiring the rule body be frequent, it requires the entire rule be frequent. To simply achieve this goal, we also need to require the line 06 of SUBROUTINE 1.1 satisfy $sup_{p_i}^{c_j} \geq min_sup$ (or $sup_{p_i}^{c_j} \geq min_sup_j$ in the case of applying varying support threshold).

Adapting the algorithm to the traditional definition of a classification rule also enables us to design some

search space pruning methods. Let the current prefix be P , a local item of P , x_j , is called infrequent and can be safely pruned according to Lemma ??.

Lemma 4 (Infrequent item pruning) *Item x_j is called an infrequent item w.r.t. prefix P and can be safely pruned from P 's conditional database if it satisfies the following equation:*

$$\max_{\forall i, 1 \leq i \leq k} sup_{P \cup \{x_j\}}^{c_i} < min_sup \quad (8)$$

Proof. Follows easily from the traditional definition of a classification rule. \square

5.4.4 Maximum Support Threshold

Some dense databases contain some highly frequent items, which appear in almost all the training instances. From the classification point of view, these items are indifferentiable and cannot be used to generate high quality classification rules. Removing these items usually does not hurt the classification accuracy, but it can significantly improve the algorithm efficiency. Thus, in HARMONY there is an option for the user to specify a maximum support threshold, max_sup , in order to remove the overly frequent items.

6 Empirical Results

6.1 Test Environment and Databases

We implemented the HARMONY algorithm in C and performed a thorough experimental study. We first evaluated HARMONY as a frequent itemset mining algorithm to show the effectiveness of the pruning methods, the algorithm efficiency and scalability. Then we compared HARMONY with some well-known classifiers on both categorical and text databases. All the experiments were performed on a 1.8GHz Linux machine with 1GB memory.

The UCI Databases. Many previous studies used some small databases to evaluate both the accuracy and efficiency of a classifier. For example, most of the 26 databases used in [?, ?, ?] only contain several hundred instances, which means the test databases contain too few test instances (i.e., only a few tens) if the 10-fold cross validation is adopted to evaluate the classification accuracy. In this paper, we mainly focus on relatively large databases (by large, we mean the database should contain no fewer than 1000 instances), although we also report the comparison results for some small databases.

Table 2. Large UCI database characteristics.

| Database | # instances | # items | # classes |
|------------|-------------|---------|-----------|
| adult | 48842 | 131 | 2 |
| chess | 28056 | 66 | 18 |
| connect | 67557 | 66 | 3 |
| led7 | 3200 | 24 | 10 |
| letRegcog | 20000 | 106 | 26 |
| mushroom | 8124 | 127 | 2 |
| nursery | 12960 | 32 | 5 |
| pageBlocks | 5473 | 55 | 5 |
| penDigits | 10992 | 90 | 10 |
| waveform | 5000 | 108 | 3 |

Table 3. Small UCI database characteristics.

| Database | # instances | # items | # classes |
|-------------|-------------|---------|-----------|
| anneal | 798 | 106 | 6 |
| auto | 205 | 142 | 7 |
| breast | 699 | 48 | 2 |
| glass | 214 | 52 | 7 |
| heart | 303 | 53 | 5 |
| hepatitis | 155 | 58 | 2 |
| horseColic | 368 | 94 | 2 |
| ionosphere | 351 | 104/173 | 2 |
| iris | 150 | 23 | 3 |
| pimaIndians | 768 | 42 | 2 |
| ticTacToe | 958 | 29 | 2 |
| wine | 178 | 68 | 3 |
| zoo | 101 | 43 | 7 |

In [?], the author used 23 UCI databases to compare FOIL and CPAR algorithms⁴. Among these 23 databases, 10 of them are large databases and the left 13 databases are small ones. The characteristics of these two classes of databases are summarized in Table ?? and Table ??, respectively. All the 23 databases were obtained from the author of [?] and the 10-fold cross validation is used for comparison with FOIL and CPAR. Among these databases, database *ionosphere* was discretized into two versions, one with 104 distinct items (denoted by *ionosphere104*) and another with 173 distinct items (denoted by *ionosphere173*). Because databases *connect* and *ionosphere104* are too dense, during the 10-fold cross validation in our experiments HARMONY only used the items whose supports are no greater than 20,000 and 190 for *connect* and *ionosphere104* respectively, to generate classification rules (i.e., $max_sup=20,000$ and $max_sup=190$ for

⁴The numerical attributes in these databases have been discretized by the author, and the discretization technique is different from those used in [?, ?, ?]; thus, the performance reported here may be different from the previous studies even for the same algorithm and the same database.

these two databases respectively).

Table 4. Top 10 topics in reuters-21578.

| Category Name | # train labels | # test labels |
|---------------|----------------|---------------|
| acq | 1650 | 719 |
| corn | 181 | 56 |
| crude | 389 | 189 |
| earn | 2877 | 1087 |
| grain | 433 | 149 |
| interest | 347 | 131 |
| money-fx | 538 | 179 |
| ship | 197 | 89 |
| trade | 369 | 118 |
| wheat | 212 | 71 |
| total | 7193 | 2787 |

Table 5. Class distribution in *sports* database.

| Class Name | Number of labels |
|------------|------------------|
| baseball | 3412 |
| basketball | 1410 |
| football | 2346 |
| hockey | 809 |
| boxing | 122 |
| bicycle | 145 |
| golf | 336 |
| total | 8580 |

Text Databases. We also used two text databases in our empirical evaluation. The first database is the popularly used ‘ModeApte’ split version of the reuters-21578 collection, which was preprocessed and provided by the authors of [?], and both the database and its description are available at [?]. After preprocessing, it contains totally 8575 distinct terms, 9603 training documents, and 3299 test documents. Like many other studies [?, ?, ?, ?], we are more interested in the top 10 most common categories (i.e., topics). These ten largest categories form 6488 training documents and 2545 test documents. A small portion of the training and test documents are associated with multiple category labels (that is, *reuter-21578* is a multi-class multi-label database). In our experiments, we treated each one of the training documents with multiple labels as multiple documents, each one with a distinct label. The top 10 categories and their corresponding number of labels in the training and test databases are described in Table ?. The second text database is *sports*, which was obtained from San Jose Mercury (TREC). In our experiments, we removed some highly frequent terms, and finally it contains totally 8580 doc-

uments, 7 classes, and about 1748 distinct terms. The seven classes and their corresponding number of documents are shown in Table ??.

6.2 Experimental Results

6.2.1 Evaluate HARMONY as a Frequent Itemset Mining Algorithm

To mine the highest confidence covering rule(s) for each instance, a naïve method is like the association-based classifiers: first use an efficient association rule mining algorithm to compute the complete set of classification rules, from which the set of the highest confidence covering rules w.r.t. each instance can be selected. Our empirical results show that this method is usually inefficient if the database is large and a more efficient way is to push some effective pruning methods into the frequent itemset mining framework and to directly mine the final set of classification rules.

Effectiveness of the pruning methods. We first evaluated the effectiveness of the pruning methods. Figure ??a shows the results for database *penDigits* with absolute support threshold varying from 512 to 8. At first glance of Equation ?? and Equation ??, the *unpromising item* and *conditional database* pruning methods seem to be less effective at lower support, however this is not the case when considering more covering rules with higher confidence can be found at lower support and can be used to more quickly raise the currently maintained highest confidences. As we can see from Figure ??a, if we turn off the pruning methods used in HARMONY (denoted by ‘without pruning’), it can become over an order of magnitude slower at low support.

Scalability test. Figure ??b shows the scalability test result for databases *letRecog*, *waveform*, and *mushroom* with relative support set at 0.5%. In the experiments, we replicated the instances from 2 to 16 times. We can see that HARMONY has linear scalability in the runtime with increasing number of instances.

Efficiency test. As we mentioned above, the traditional frequent (closed) itemset mining algorithms can be revised to mine the complete set of high confidence classification rules, from which a subset of high quality rules can be further identified. Our efficiency tests for HARMONY in comparison with FPgrowth* and FPclose, two recently developed efficient frequent/closed itemset mining algorithms [?], show that such a method is not realistic at low support, while our experiments demonstrate that the classification accuracy is usually higher at lower support.

Figure ?? shows the comparison results for database *sports*. As we can see, although at high support, both

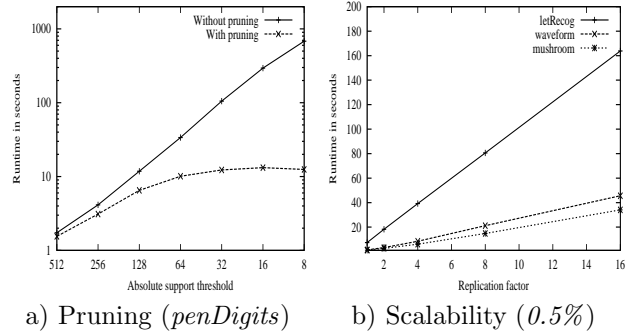


Figure 1. Pruning and scalability test.

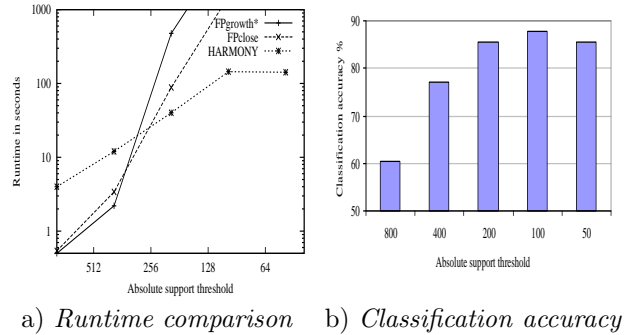


Figure 2. Efficiency test (*sports*).

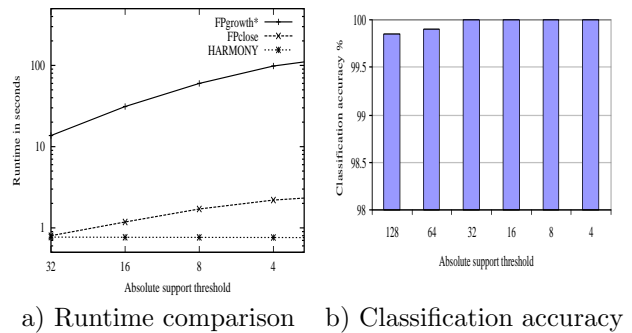


Figure 3. Efficiency test (*mushroom*).

FPgrowth* and FPclose are faster than HARMONY, once we continue to lower the support, they will be much slower. For example, at absolute support of 100, HARMONY is several orders of magnitude faster than FPgrowth* and FPclose. Figure ??b shows the classification accuracy at different support thresholds using the 10-fold cross validation. We can see that HARMONY can achieve higher accuracy at lower support like 100. It is also interesting to see that the accuracy at a too low support 50 is worse than that at support 100 for this database, due to the ‘overfitting’ problem.

Figure ??a shows similar comparison results for categorical database *mushroom*. HARMONY is faster

than both FPgrowth* and FPclose at absolute support lower than 32. Figure ??b shows that HARMONY has better accuracy at lower support threshold.

6.2.2 Classification Evaluation

The reuters-21578 (ModApte) text database.

For a multi-class multi-label database like *reuters-21578*, most previous studies used the breakeven point of precision and recall to measure the classifier performance [?, ?, ?, ?, ?, ?, ?], which is defined as the point at which precision is equal to the recall. To our best knowledge, the best breakeven performance for the *reuters-21578* database is the linear SVM [?]. For comparison with earlier results, we first found the overall breakeven point in terms of all the top 10 categories by adjusting the dominant factor γ , then reported the average of precision and recall for each category as their corresponding breakeven performance [?].

Table ?? shows the comparison results with some previous results. The results for Findsim (i.e., Find-Similar), NBayes (i.e., Naïve-Bayes), Bayes-Nets, Trees (i.e., Decision-Trees), and LinearSVM were obtained from [?]. The micro-avg is the overall breakeven performance over all 10 categories. For HARMONY, we used three different uniform absolute support thresholds, 60, 70, and 80, respectively. From Table ?? we can see that both HARMONY and LinearSVM have similar breakeven performance and perform much better than all the other classifiers, including Find-Similar, Naïve-Bayes, Bayes-Nets, and Decision-Trees. Among the 10 categories, HARMONY achieves the best performance at support of 60 for five categories, *acq*, *earn*, *money-fx*, *ship*, and *trade*. While LinearSVM performs best for another three categories, *crude*, *grain*, and *interest*. Decision-Trees also performs good and has the best performance for two small categories, *corn* and *wheat*. SVM is very well known for classifying high dimensional text databases. Our results show that HARMONY can achieve similar performance to SVM. Table ?? shows the runtime in seconds for HARMONY at three different support thresholds. We can see that HARMONY is very efficient in classifying the *reuters-21578* text database. For example, at absolute support of 60, it takes 72.6 seconds to build the model from 7193 training documents, and 0.363 seconds to classify 2545 test documents.

The micro-avg values for HARMONY in Table ?? were computed by using a uniform minimum support for all 10 categories. If we choose some proper varying support thresholds for different categories, HARMONY can achieve better performance. In Table ??, the second row shows the corresponding minimum sup-

port chosen for the 10 categories, while the third row shows the breakeven performance. From the results we can see that with these varying support thresholds, HARMONY achieves a better micro-avg value, 92.4. This example illustrates that adopting different support thresholds for different classes does achieve better results.

The micro-avg values for HARMONY in Table ?? and Table ?? were computed by setting the *score differentia factor* at its default value 1. By choosing different *differentia factor* values, HARMONY may have different micro-avg performance. Table ?? shows the micro-avg performance and the corresponding dominant factor γ for HARMONY with the varying support thresholds shown in Table ?? and by varying the parameter of *score differentia factor* δ from 0 to 1.2. $\delta = 0$ means the weighted dominant factor-based scoring method degenerates to the dominant factor-based method used in [?]. By adopting a proper value of δ , the weighted dominant factor-based scoring method can achieve a better micro-avg performance. For example, by setting δ at 0.9, the overall precision equals the overall recall at $\gamma = 0.544$, and the corresponding micro-avg breakeven performance for HARMONY is 92.4, which is higher than the corresponding micro-avg at $\delta = 0$ (i.e., 91.2).

The UCI databases. We evaluated HARMONY on the UCI databases in comparison with FOIL and CPAR, which are two well-known algorithms for classifying categorical data. The results in [?] show that CPAR has comparable accuracy to the association-based algorithms CMAR and CBA, but is more efficient; thus, we will do not compare HARMONY with the association-based algorithms. The results for FOIL and CPAR were provided by Frans Coenen and are available at [?]. All the results including the accuracy, runtime, and the number of rules are computed using the 10-fold cross validation. The reported accuracy and number of rules are the corresponding average value of the 10-fold cross validation results, while the runtime is the total runtime of the 10-fold cross validation, including both training and testing time. In the experiments, we fixed the absolute support threshold at 50 for HARMONY with all 10 large UCI databases, and at 10 for all 13 small UCI databases.

Table ?? shows the comparison results in terms of the 10-fold cross validation accuracy for 10 large UCI databases. These results show that HARMONY has better accuracy than both FOIL and CPAR for most of 10 large UCI databases, and has comparable accuracy with FOIL for databases *adult*, *mushroom*, and *pageBlocks*. On average, HARMONY has significantly better accuracy than both FOIL and CPAR: its average

Table 6. Breakeven performance on the *Reuters-21578* database with some well-known classifiers.

| Categories | HARMONY <i>min_sup=60</i> | HARMONY <i>min_sup=70</i> | HARMONY <i>min_sup=80</i> | Findsim | NBayes | BayesNets | Trees | SVM (linear) |
|------------|------------------------------|------------------------------|------------------------------|---------|--------|-----------|-------------|-----------------|
| acq | 95.3 | 95.3 | 95.3 | 64.7 | 87.8 | 88.3 | 89.7 | 93.6 |
| corn | 78.2 | 78.6 | 75.2 | 48.2 | 65.3 | 76.4 | 91.8 | 90.3 |
| crude | 85.7 | 85.0 | 88.0 | 70.1 | 79.5 | 79.6 | 85.0 | 88.9 |
| earn | 98.1 | 98.2 | 97.6 | 92.9 | 95.9 | 95.8 | 97.8 | 98.0 |
| grain | 91.8 | 90.4 | 90.1 | 67.5 | 78.8 | 81.4 | 85.0 | 94.6 |
| interest | 77.3 | 76.6 | 75.1 | 63.4 | 64.9 | 71.3 | 67.1 | 77.7 |
| money-fx | 80.5 | 81.9 | 82.1 | 46.7 | 56.6 | 58.8 | 66.2 | 74.5 |
| ship | 86.9 | 82.9 | 82.8 | 49.2 | 85.4 | 84.4 | 74.2 | 85.6 |
| trade | 88.4 | 88.0 | 86.1 | 65.1 | 63.9 | 69.0 | 72.5 | 75.9 |
| wheat | 62.8 | 60.6 | 58.7 | 68.9 | 69.7 | 82.7 | 92.5 | 91.8 |
| micro-avg | 92.0 | 91.7 | 91.4 | 64.6 | 81.5 | 85.0 | 88.4 | 92.0 |

Table 7. Breakeven performance on the *Reuters-21578* database with varying *min_sup*.

| Categories | acq | corn | crude | earn | grain | interest | money-fx | ship | trade | wheat | micro-avg |
|----------------|------|------|-------|------|-------|----------|----------|------|-------|-------|-----------|
| <i>min_sup</i> | 55 | 60 | 60 | 75 | 60 | 55 | 70 | 50 | 60 | 45 | - |
| Breakeven | 95.6 | 78.2 | 86.6 | 97.9 | 90.7 | 76.8 | 83.5 | 88.5 | 89.3 | 68.1 | 92.4 |

accuracy over all 10 large UCI databases is about 5% higher than FOIL, and 10% higher than CPAR. Note in the experiments we fixed the minimum support at 50 for all 10 large UCI databases, if we choose a lower (or higher) supports for some databases, HARMONY can achieve much better performance than what we reported here. For example, if we choose the minimum support at 5 for *chess* database, HARMONY has a classification accuracy 58.43%, which is more than 13% higher than the accuracy at support 50, while it only becomes about two times slower.

Table 8. Runtime on *reuters-21578* database.

| Runtime | HARMONY (<i>min_sup=60</i>) | HARMONY (<i>min_sup=70</i>) | HARMONY (<i>min_sup=80</i>) |
|----------|----------------------------------|----------------------------------|----------------------------------|
| training | 72.6 | 51.1 | 37.6 |
| testing | 0.363 | 0.337 | 0.309 |

Table 9. Effectiveness of the *score differentia* factor δ (varying *min_sup*, *reuters-21578*).

| δ | 0 | 0.3 | 0.6 | 0.9 | 1.2 |
|-----------|--------|--------|-------|-------|-------|
| γ | 0.4987 | 0.5029 | 0.517 | 0.544 | 0.569 |
| micro-avg | 91.2 | 91.6 | 92.2 | 92.4 | 92.1 |

Table ?? compares the runtime (in seconds) of the three algorithms on 10 large UCI databases. Note that FOIL and CPAR were implemented in java and were tested on a different machine from that of HARMONY. As a result, these times cannot be directly compared to the times reported for HARMONY but they only provide an overall idea on the relative computational requirements of the various schemes. Table ?? shows

Table 10. Accuracy comparison on 10 large UCI databases (*min_sup=50* for HARMONY).

| Database | FOIL | CPAR | HARMONY |
|------------|-------------|-------|---------------|
| adult | 82.5 | 76.7 | 81.9 |
| chess | 42.6 | 32.8 | 44.87 |
| connect | 65.7 | 54.3 | 68.05 |
| led7 | 62.3 | 71.2 | 74.56 |
| letRecog | 57.5 | 59.9 | 76.81 |
| mushroom | 99.5 | 98.8 | 99.94 |
| nursery | 91.3 | 78.5 | 92.83 |
| pageBlocks | 91.6 | 76.2 | 91.6 |
| penDigits | 88.0 | 83.0 | 96.23 |
| waveform | 75.6 | 75.4 | 80.46 |
| average | 75.66 | 70.68 | 80.725 |

that on average the runtime of HARMONY is over an order of magnitude smaller than those of FOIL and CPAR. For some large databases like *chess*, the runtime of HARMONY can be over two orders of magnitude smaller than those of FOIL and CPAR. Table ?? compares the number of rules discovered by three algorithms. We can see that on average, HARMONY finds many more rules than both FOIL and CPAR. The reason why HARMONY finds more rules is that it mines classification rules in an instance-centric manner: it mines at least one highest confidence covering rule for each instance.

Table ?? depicts the accuracy comparison among FOIL, CPAR, and HARMONY on 13 small UCI databases, from which we can see that on average HARMONY and FOIL have similar classification accuracy and both perform a little better than CPAR. In

Table 11. Runtime comparison on 10 large UCI databases ($min_sup=50$ for HARMONY).

| Database | FOIL | CPAR | HARMONY |
|------------|---------|--------------|----------------|
| adult | 10251.0 | 809.0 | 1395.5 |
| chess | 10122.8 | 1736.0 | 11.34 |
| connect | 35572.5 | 24047.1 | 85.44 |
| led7 | 11.5 | 5.7 | 1.29 |
| letRecog | 4365.6 | 764.0 | 778.91 |
| mushroom | 38.3 | 15.4 | 8.78 |
| nursery | 73.1 | 51.7 | 6.21 |
| pageBlocks | 43.1 | 15.5 | 2.5 |
| penDigits | 821.1 | 101.9 | 82.6 |
| waveform | 295.3 | 38.1 | 130.0 |
| total | 61594.3 | 27584.4 | 2502.57 |

Table 12. Comparison of # rules on 10 large UCI databases ($min_sup=50$ for HARMONY).

| Database | FOIL | CPAR | HARMONY |
|------------|--------|--------|---------|
| adult | 331.8 | 183.1 | 6431.3 |
| chess | 1116.7 | 1504.8 | 2881.5 |
| connect | 285.8 | 816.1 | 6664.0 |
| led7 | 80.6 | 31.4 | 268.7 |
| letRecog | 560.9 | 643.0 | 2255.8 |
| mushroom | 16.2 | 30.8 | 95.9 |
| nursery | 57.4 | 83.6 | 391.64 |
| pageBlocks | 123.1 | 56.2 | 78.1 |
| penDigits | 204.6 | 166.9 | 1434.5 |
| waveform | 159.7 | 114.3 | 958.6 |
| average | 293.68 | 363.02 | 2146.0 |

the experiments we fixed the minimum absolute support at 10 on all 13 small UCI databases for HARMONY, if we use tuned minimum support, HARMONY can achieve better accuracy for most databases. In addition, because these databases contain a small number of training instances, the number of classification rules mined by HARMONY may not be sufficient to build an accurate classification model; thus, we implemented a variant of HARMONY, which mines K highest confidence frequent covering rules for each training instance if it supports no fewer than K such rules. By varying K parameter from 1 to 5, and choosing the minimum absolute support from $\{5, 10, 15\}$, we got a set of classification results, among which the best results for 13 small UCI databases are shown in Table ???. We can see that with tuned parameters, HARMONY achieves better classification accuracy than both FOIL and CPAR.

We also evaluated the effectiveness of the *support differentia factor* ξ using one unbalanced UCI database, *hepatitus*. In Table ??, we set min_sup at 10 and varied ξ from 0 to 1. We can see that the *support differentia factor* based varying support threshold method is very effective in improving the accuracy for this database. For example, with $\xi = 1$, HARMONY has an accuracy 84.84, which is much higher than that at $\xi = 0$.

Table 13. Accuracy comparison on 13 small UCI databases ($min_sup=10$ for HARMONY).

| Database | FOIL | CPAR | HARMONY |
|-------------|-------------|-------------|--------------------|
| anneal | 96.9 | 90.2 | 91.51 |
| auto | 46.1 | 48.0 | 61 |
| breast | 94.4 | 94.8 | 92.42 |
| glass | 49.3 | 48.0 | 49.8 |
| heart | 57.4 | 51.1 | 56.46 |
| hepatitus | 77.5 | 76.5 | 78 |
| horseColic | 83.5 | 82.3 | 82.53 |
| ionosphere | 89.5 | 92.9 | 92.03/88.31 |
| iris | 94.0 | 94.7 | 93.32 |
| pimaIndians | 73.8 | 75.6 | 72.34 |
| ticTacToe | 96.0 | 72.2 | 92.29 |
| wine | 86.4 | 92.5 | 91.94 |
| zoo | 96.0 | 96.0 | 93.0 |
| average | 80.06 | 78.06 | 80.43/80.14 |

7 Conclusion

Designing accurate, efficient, and scalable classifiers is an important research topic in data mining. In many applications, the databases are high-dimensional. Due to the curse of dimensionality, some traditional classification algorithms may not work well for this type of data. The rule-based classifiers provide a promising approach to tackle this problem. However, to achieve high accuracy, a good rule-based classifier needs to find a sufficient number of high quality classification rules and use them to build the model.

In this paper, we proposed an instance-centric classification rule mining paradigm and designed an accurate classifier, HARMONY. Several effective search space pruning methods have also been proposed, which can be pushed deeply into the projection-based frequent itemset enumeration framework. Our performance study shows that HARMONY has high accuracy and efficiency in comparison with many well known classifiers for both categorical data and high dimensional text data. It also has good scalability in terms

Table 14. Tuned accuracy on 13 small UCI databases for HARMONY.

| Database | min_sup $\in \{5, 10, 15\}$ | K rules $\in [1, 5]$ | HARMONY |
|-------------------------|-----------------------------------|-------------------------|--------------------|
| anneal | 5 | 4 | 95.65 |
| auto | 10 | 1 | 61.5 |
| breast | 15 | 2 | 96.14 |
| glass | 10 | 1 | 49.8 |
| heart | 15 | 1 | 58.4 |
| hepatitus | 15 | 5 | 84.41 |
| horseColic | 5 | 4 | 84.64 |
| ionosphere (104/173) | 10/15 | 2/3 | 93.45/88.85 |
| iris | 5 | 5 | 95.99 |
| pimaIndians | 5 | 5 | 73.79 |
| ticTacToe | 10 | 4 | 94.09 |
| wine | 10 | 5 | 94.9 |
| zoo | 5 | 1 | 96 |
| average | | | 82.98/82.62 |

Table 15. Effectiveness of the support differentia factor ξ ($min_sup=10$, hepatitus).

| ξ | 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1 |
|----------|----|-------|-----|------|-----|-------|
| accuracy | 78 | 83.34 | 84 | 81.5 | 83 | 84.84 |

of database size.

Acknowledgements

We are grateful to Frans Coenen at the University of Liverpool and Shane Bergsma at the University of Alberta for providing us the discretized UCI databases and the reuters-21578 database, respectively, and promptly answering our various questions. We also thank Osmar R. Zaiane and Maria-Luiza Antonie at the University of Alberta for answering our questions related to the ARC-BC algorithm.

References

- [1] R. Agarwal, C. Aggarwal, V. Prasad. *A Tree Projection Algorithm for Generation of Frequent Item Sets*, Journal of Parallel and Distributed Computing. 61(3), 2001.
- [2] R. Agrawal, T. Imielinski, A. Swami. *Mining Association Rules between Sets of Items in Large Databases*, SIGMOD'93.
- [3] R. Agrawal, R. Srikant. *Fast Algorithms for Mining Association Rules*, VLDB'94.
- [4] K. Ali, S Manganaris, R. Srikant. *Partial Classification Using Association Rules*, KDD'97.
- [5] M. Antonie, O. Zaiane. *Text Document Categorization by Term Association*, ICDM'02.
- [6] C. Apte, F. Damerau, S.M. Weiss. *Towards Language Independent Automated Learning of Text Categorization Models*, SIGIR'94.
- [7] R.J. Bayardo. *Brute-force Mining of High-confidence Classification rules*, KDD'97.
- [8] R.J. Bayardo, R. Agrawal. *Mining the most interesting rules*, KDD'99.
- [9] R. Bekkerman, R. El-Yaniv, N. Tishby, Y. Winter. *On Feature Distribution Clustering for Text Categorization*, SIGIR'01.
- [10] S. Bergsma. *The Reuters-21578 (ModApte) dataset*, Department of Computer Science, University of Alberta. Available at <http://www.cs.ualberta.ca/~bergsma/650/>.
- [11] S. Bergsma, D. Lin. *Title Similarity-Based Feature Weighting for Text Categorization*, CMPUT 650 Research Project Report, Department of Computer Science, University of Alberta.
- [12] F. Coenen. (2004) The LUCS-KDD Implementations of the FOIL, PRM, and CPAR algorithms, http://www.csc.liv.ac.uk/~frans/KDD/Software/FOIL_PRM_CPAR/foilPrmCpar.html, Computer Science Department, University of Liverpool, UK.
- [13] W. Cohen. *Fast effective rule induction*, ICML'95.
- [14] G. Cong, X. Xu, F. Pan, A. Tung, J. Yang. *FARMER: Finding Interesting Rule Groups in Microarray Datasets*, SIGMOD'04.
- [15] M. Deshpande, G. Karypis. *Using Conjunction of Attribute Values for Classification*, CIKM'02.
- [16] S. Dumais, J. Platt, D. Heckerman, M. Sahami. *Inductive Learning Algorithms and Representations for Text Categorization*, CIKM'98.
- [17] T. Fukuda, Y. Morimoto, S. Motishita. *Constructing Efficient Decision Trees by Using Optimized Numeric Association Rules*, VLDB'96.
- [18] K. Gade, J. Wang, G. Karypis. *Efficient Closed Pattern Mining in the Presence of Tough Block Constraints*, KDD'04.
- [19] G. Grahne, J. Zhu. *Efficiently Using Prefix-trees in Mining Frequent Itemsets*, ICDM-FIMI'03.
- [20] J. Han, J. Pei, Y. Yin. *Mining Frequent Patterns without Candidate Generation*, SIGMOD'00.

- [21] T. Joachims. *Text Categorization with Support Vector Machines: Learning with Many Relevant Features*, ECML'98.
- [22] B. Lent, A. Swami, J. Widom. *Clustering Association Rules*, ICDE'97.
- [23] N. Lesh, M. Zaki, M. Ogihara. *Mining Features for Sequence Classification*, KDD'99.
- [24] W. Li, J. Han, J. Pei. *CMAR: Accurate and Efficient Classification based on multiple class-association rules*, ICDM'01.
- [25] B. Liu, W. Hsu, Y. Ma. *Integrating Classification and Association Rule Mining*, KDD'98.
- [26] J. Quinlan. *C4.5: Programs for Machine Learning*, Morgan Kaufman, 1993.
- [27] J. Quinlan, R. Cameron-Jones. *FOIL: A Midterm Report*, ECML'93.
- [28] J. Wang, G. Karypis. *BAMBOO: Accelerating Closed Itemset Mining by Deeply Pushing the Length-Decreasing Support Constraint*, SDM'04.
- [29] Y. Yang. *An Evaluation of Statistical Approaches to Text Categorization*, Information Retrieval, Vol. 1, No. 1-2, 1999.
- [30] X. Yin, J. Han. *CPAR: Classification based on Predictive Association Rules*, SDM'03.
- [31] M. Zaki, C. Aggarwal. *XRULES: An Effective Structural Classifier for XML Data*, KDD'03.
- [32] Y. Zhao, G. Karypis. *Criterion Functions for Document Clustering: Experiments and Analysis*, Machine Learning, in press.