

On Optimal Rule Discovery

Jiuyong Li

Department of Mathematics and Computing

The University of Southern Queensland

Australia, 4350

Jiuyong.Li@usq.edu.au

September 25, 2005

Abstract

In machine learning and data mining, heuristic and association rules are two dominant schemes for rule discovery. Heuristic rule discovery usually produces a small set of accurate rules, but fails to find many globally optimal rules. Association rule discovery generates all rules satisfying some constraints, but yields too many rules and is infeasible when the minimum support is small. Here we present a unified framework for the discovery of a family of optimal rule sets, and characterise the relationships with other rule discovery schemes such as non-redundant association rule discovery. We theoretically and empirically show that optimal rule discovery is significantly more efficient than the association rule discovery independent of data structure and implementation. Optimal rule discovery is an efficient alternative to association rule discovery, especially when the minimum support is low.

keywords: Data mining, rule discovery, optimal rule set.

1 Introduction

Rules are one of the most expressive and human understandable representations of knowledge; a rule based method produces self-explanatory results. Therefore, rule discovery has been a major issue in machine learning and data mining.

Heuristic algorithms for rule discovery that were developed in the machine learning community, such as C4.5rules [15], CN2 [6], and RIPPER [7] focus on classification accuracy and usually return small rule sets. However, a heuristic method does not guarantee the discovery of the best quality rules. A complete or optimal rule set is more desirable whenever it is computationally feasible.

Association rule discovery [1] produces a complete rule set within the minimum support and confidence constraints. It has been widely accepted because of the simplicity of the problem statement and the effectiveness of pruning by support. Association rule discovery is a general purpose rule discovery scheme, and has wide applications. Classification is one application. CBA [12] makes use of a method that is similar to the C4.5rules pruning method to prune an association rule set and produces more accurate classifiers than C4.5rules. This suggests that some rules in the complete rule set, which are missed by C4.5rules, make CBA classifiers more accurate. However, association rule discovery usually produces too many rules and is inefficient when the minimum support is low.

Non-redundant association rule discovery [19] improves the efficiency of association rule discovery. However, the requirements for redundant rules are strict, and the efficiency of non-redundant rule discovery can be further improved. We discuss the relationships between our proposed optimal rule set and the redundant rule set in Section 3.

Optimal rule discovery uncovers rules that maximise an interestingness measure. The search for maxima further prunes the search space, and hence optimal rule discov-

ery is significantly more efficient than association rule discovery.

One type of optimal rule sets is k largest rule sets, which contain the top k rules measured by an interestingness metric. Webb and Zhang's k -optimal rule set [18] is a typical example. k -optimal rules are measured by a Leverage metric. However, the top k rules may come from the same section of data, and leave some records in a data set uncovered by rules. This is a drawback for optimal rule sets containing k largest rules.

The problem of not reasonably covering the data set exists in other optimal rule sets, such as the SC optimality rule set [2] and rule sets defined by all confidence and bond [13]. In a SC optimality rule set, a rule with higher confidence and support excludes another rule with lower confidence and support. When the records covered by the excluded rule are not covered by another rule, these records lose their representative rules in the SC-optimal rule set. In the generation of rule sets defined by all confidence and bond, rules with different targets are compared directly for the exclusion of rules from the rule sets. Rules with different targets have different implications and they should not be used to exclude each others.

The definition of optimal rule sets in this paper is very close to a special constraint rule set with a zero confidence improvement [3], which consists of rules whose confidences are greater than confidences of all their simpler form rules. A rule covers a subset of records covered by one of its simpler form rule, and hence it is guaranteed that the records covered by the excluded rule are covered by other rules with higher confidence. A PC optimality rule set [2] post prunes the constraint association rule set with a zero confidence improvement [3], and hence is in the same category of the optimal rule sets which we discuss.

We achieve the following four developments in this paper. First, we present a general definition for a family of optimal rule sets for a range of interestingness metrics. Second, we prove that the family of optimal rule sets observe the same anti-monotonic property. Third, we develop an effective algorithm for mining the family of optimal

rule sets without assuming that the target is fixed to one class. Fourth, we characterise the relationships between an optimal rule set and a non-redundant rule set, and reveal the relationships among support pruning, closure pruning, and optimality pruning.

The rest of this paper is arranged as follows: Section 2 presents definitions; Section 3 gives properties of the family of optimal rule sets; Section 4 provides a complete algorithm for mining the optimal rule sets; Section 5 illustrates the relationships among support pruning, closure pruning and optimality pruning; Section 6 presents proof-of-concept experimental results; and Section 7 concludes the paper.

2 Definitions

Consider a relational dataset D with n attribute domains. A record of D is set of attribute-value pairs, denoted by T . A *pattern* is a subset of a record. We say a pattern is a k -*pattern* if it contains k attribute-value pairs. All the records in D are categorised by a set of classes C .

An *implication* is denoted by $P \rightarrow c$, where P is called the *antecedent*, and c is called the *consequence*. The *support* of pattern P is defined to be the ratio of the number of records containing P to the number of all the records in D , denoted by $\text{supp}(P)$. The support of implication $P \rightarrow c$ is defined to be the ratio of the number of records containing both P and c to the number of all the records in D , denoted by $\text{supp}(P \rightarrow c)$. The *confidence* of the implication $P \rightarrow c$ is defined to be the ratio of $\text{supp}(P \rightarrow c)$ to $\text{supp}(P)$, represented by $\text{conf}(P \rightarrow c)$.

An *association rule* is a strong implication whose both support and confidence are not less than given thresholds from a dataset.

The cover set of pattern P is the set of IDs of records containing P , represented by $\text{cov}(P)$. The cover set of rule $P \rightarrow c$ is a set of IDs of records containing both P and c , denoted by $\text{cov}(P \rightarrow c)$. Clearly, if $P \subset Q$ then we have $\text{cov}(P) \supseteq \text{cov}(Q)$

and $\text{cov}(P \rightarrow c) \supseteq \text{cov}(Q \rightarrow c)$.

For simplicity, in the rest of this paper we use upper case letters, for example, P and Q , to stand for patterns, and lower case letters, for example, a, b, \dots , to stand for attribute-value pairs. We abbreviate $P \cup Q$ as PQ and $P \cup \{a\}$ as Pa .

The association rule definition is understandable, but it has the following major obstacles in real world applications:

1. The confidence is not suitable for a variety of applications;
2. The number of association rules is too many; and
3. The support pruning is not efficient when the minimum support is low.

To overcome the first obstacle, many interestingness metrics have been proposed to measure interestingness of rules. For example, lift (interest or strength), gain, added-value, Klogen, conviction, p-s, Laplace, cosine, certainty factor, Jaccard, and many others discussed by Tan et al [16].

All these interestingness metrics are used monotonically. A rule with a higher value in a metric is more interesting than a rule with a lower one. To generalise, we use *Interest* to stand for an interestingness metric and $\text{Interestingness}(P \rightarrow c)$ for the interestingness of rule $P \rightarrow c$.

Some interestingness metrics are not monotonic with the interestingness, such as, odds ratio. A rule with the odds ratio that is significantly greater than or less than 1 is interesting. For example, let c be a disease and A be a symptom or exposure factor. A high odds ratio of rule $A \rightarrow c$ means the disease has higher occurrence probability in the cohort with A than the cohort without A , and vice versa. In this paper, we consider an odds ratio that is greater than 1, and hence it is monotonic with the interestingness. When we need rules with odds ratios lower than 1, we switch the consequence. For

example, when we divide data into the disease group and the non-disease group. A high odds ratio in the non-disease group means a low odds ratio in the disease group.

To make the rule definition more general, we replace the confidence by the value of an interest metric. Formally, we have the following definition.

Definition 1 *The general rule*

A rule is a strong implication whose both support and interestingness are not less than given thresholds.

In the rest of this paper, a rule refers to a generalised rule instead of an association rule. To proceed the discussions of obstacles 2 and 3, we give another definition.

Definition 2 *General and specific relationships*

Given two rules $P \rightarrow c$ and $Q \rightarrow c$ where $P \subset Q$, we say that the latter is more specific than the former and the former is more general than the latter.

A specific rule covers a subset (at most the equal set) of records covered by one of its more general rules. More formally, $\text{cov}(Q \rightarrow c) \subseteq \text{cov}(P \rightarrow c)$. Therefore, the removal of a specific rule from a rule set does not reduce the total coverage of the rule set.

In some cases, we may consider the rule $\emptyset \rightarrow c$ as the most general rule targeting c . Its confidence equals $\text{supp}(c)$. For example, if 80% customers buy bread when they shop in a supermarket, then this is formalised as $\emptyset \rightarrow \text{bread}$ (confidence = 80%). Such a rule filters many trivial rules that do not surprise a manager, for example rule $\text{egg} \rightarrow \text{bread}$ (confidence = 75%).

In other cases, we need to consider 1-pattern antecedent rules, such as $a \rightarrow c$ and $b \rightarrow c$, as the most general rules. For example, 99.9% of records in a medical data set are not related to a particular disease, but we are still interested in rules with 80% confidence in the records since they may reveal some preventative patterns. In this

paper, we consider this case.

Obstacles 2 and 3 are closely related. A major reason for a lot of rules being of no interest to users is that they are superfluous. For example, suppose that we have two rules, $(\text{salary} > \$30,000) \rightarrow \text{creditCard}(\text{approval})$ (conf = 85%) and $(\text{salary} > \$30,000)$ and $(\text{occupation} = X) \rightarrow \text{creditCard}(\text{approval})$ (conf = 84%). The latter rule is superfluous and should be removed.

There are two cases where the latter rule will be interesting: when it has much higher confidence, or when it has much lower confidence than the former rule. The primary goal for rule discovery is to find rules with the high interestingness. After identifying a small set of highly interesting rules, their exceptional rules, which have low interestingness, are considered. For example, $(\text{salary} > \$30,000)$ and $(\text{occupation} = X) \rightarrow \text{creditCard}(\text{approval})$ (conf = 30%) is an exception of rule $(\text{salary} > \$30,000) \rightarrow \text{creditCard}(\text{approval})$ (conf = 85%). After a small set of rules that are of interest to users has been identified, finding their exceptional rules is relatively simple.

Therefore, we disregard those superfluous rules in the rule discovery stage. To achieve this goal, we have the following definition.

Definition 3 *An optimal rule set*

A rule set is optimal with respect to an interestingness metric if it contains all rules except those with no greater interestingness than one of its more general rule.

Since only more specific rules are removed from an optimal rule set, an optimal rule set covers the same set of records covered by its corresponding complete rule set.

Each interestingness metric defines an optimal rule set, and the above definition defines a family of optimal rule sets. We use the following example to elaborate the definition.

Example 1 *Let interestingness metric be odds ratio (or). We have a rule set and its corresponding optimal rule set as follows.*

<i>Rule set</i>	<i>Optimal rule set</i>
$a \rightarrow z(2.3)^*, b \rightarrow z(2.0), c \rightarrow z(1.8)$ $ab \rightarrow z(2.7), ac \rightarrow z(2.1), bc \rightarrow z(1.9)$ $abc \rightarrow z(2.5)$	$a \rightarrow z(2.3), b \rightarrow z(2.0), c \rightarrow z(1.8)$ $ab \rightarrow z(2.7)$
*Numbers in parentheses are odds ratios.	

Rules $ac \rightarrow z(\text{or} = 2.1), bc \rightarrow z(\text{or} = 1.9), abc \rightarrow z(\text{or} = 2.5)$ are excluded since their odds ratios are smaller than those of their more general rules.

We provide another example to show the practical implication of an optimal rule set.

Example 2 *When the Interestingness is measured by an estimated accuracy of a rule, the optimal rule set is an optimal class association rule set [11]. Based on an ordered rule based classification model, all rules excluded by the optimal class rule set will not be used in building a classifier because they are less accurate and more complex than some rules in the optimal class association rule set covering the same data section. Therefore, a classifier built from the optimal class association rule set is identical to that from a class association rule set. In summary [11]: the optimal class association rule set is the minimum subset of rules with the same predictive power as the complete class association rule set. Further experimental proofs are reported in [10]. Classifiers built on the optimal class association rule sets are at least of the same accuracy as those from CBA [12] and C4.5rules [15].*

Building classifiers from optimal class association rule sets is significantly more efficient than from class association rule sets. Firstly, mining optimal class association rule sets is significantly faster than mining class association rule sets. When the minimum support is low, mining class association rule sets may not be feasible. Secondly, an optimal class association rule set is significantly smaller than a class association

rule set, and hence building a classifiers from the optimal class association rule set is more efficient.

In the next section, we discuss properties of the family of the optimal rule sets.

3 Properties of optimal rule sets

In this section, we discuss some properties of the family of optimal rule sets and their relationships with the non-redundant rule set.

We start with some notation. Let X be a pattern where $X \neq \emptyset$ and c be a class. PX is a proper super pattern of P . PQ is a super pattern of P . $PQ = P$ and $PQX = PX$ when $Q = \emptyset$. PQX is a proper super pattern of P . $\neg c$ stands for a special class occurring in a record where c does not occur, and therefore we have $\text{supp}(\neg c) = 1 - \text{supp}(c)$. So is it for pattern P : $\text{supp}(\neg P) = 1 - \text{supp}(P)$. $P\neg X$ is a pattern with the following support: $\text{supp}(P\neg X) = \text{supp}(P) - \text{supp}(PX)$. Further, we have $\text{supp}(\neg(PX)) = \text{supp}(\neg PX) + \text{supp}(\neg P\neg X) + \text{supp}(P\neg X)$.

Theorem 1 *Anti-monotonic property*

if $\text{supp}(PX\neg c) = \text{supp}(P\neg c)$ ¹ then rule $PX \rightarrow c$ and all its more specific rules will not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Klogsen, conviction, p-s (or leverage), Laplace, cosine, certainty factor or Jaccard.

A proof is provided in the appendix.

The practical implication of the above theorem is that: once $\text{supp}(PX\neg c) = \text{supp}(P\neg c)$ is observed, it is not necessary to search for more specific rules of $PX \rightarrow c$, for example $PQX \rightarrow c$. Those more specific rules will not be in an optimal rule set.

Rule $PX \rightarrow c$ is removed since it is not in an optimal rule set either.

¹In this paper, we only discuss rules with a single class as the consequence, but this theorem holds for rules with a pattern as the consequence.

In the following, we consider two special cases of the above theorem.

Corollary 1 *Closure property*

If $\text{supp}(P) = \text{supp}(PX)$, then rule $PX \rightarrow c$ for any c and all its more specific rules do not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Klogsen, conviction, p-s (or leverage), Laplace, cosine, certainty factor or Jaccard.

A proof is provided in the appendix.

The practical implication of the above corollary is that: once $\text{supp}(PX \neg c) = \text{supp}(P \neg c)$ is observed, it is not necessary to search for rules with PQX as their antecedent for any Q . Those rules will not be in the optimal rule set.

The reason for naming the above corollary as closure property is that it is closely related to closed pattern set mining [20].

Pattern P^c is closed if there exists no proper super pattern $X \supset P^c$ such that $\text{cov}(X) = \text{cov}(P^c)$. Consider a chain of patterns $P \subset P' \subset P'' \subset P^c$ which satisfies $\text{cov}(P) = \text{cov}(P') = \text{cov}(P'') = \text{cov}(P^c)$. P^c is the closure of patterns P , P' , and P'' . A closed pattern is the same as its closure. The support of a pattern is equivalent to that of its closure. In the above example, $\text{supp}(P) = \text{supp}(P') = \text{supp}(P'') = \text{supp}(P^c)$. Further, a pattern Y is a proper generator of Y' if $Y' \supset Y$ and $\text{cov}(Y') = \text{cov}(Y)$ hold. A pattern is a minimal generator if it has no proper generator. For example, P' is a generator of P'' and P^c , and P is a minimal generator of P^c if there exists no $Z \subset P$ such that $\text{cov}(P) = \text{cov}(Z)$.

The closed pattern and minimal generator are two closely related concepts and the mining methods for both are very similar. The condition of Corollary 1 is the fundamental test for mining both patterns. For example, we have $\text{supp}(P) = \text{supp}(PX)$ for any $X \subseteq (P^c \setminus P)$ in the above example. We illustrate this in Section 5. Usually, closed patterns are useful for producing all frequent patterns and minimal generators

are useful for generating non-redundant rules.

Zaki [19] gave a general definition of non-redundant association rule sets. Here, we rephrase it in a simple form by constraining the consequence to c . Rule $Y \rightarrow c$ is redundant if there exists $Y \supset X$ such that $\text{cov}(Y) = \text{cov}(X)$. $\text{supp}(Y) = \text{supp}(X)$ is the immediate result of $\text{cov}(Y) = \text{cov}(X)$. For example, in the chain $P \subset P' \subset P'' \subset P^c$, all rules, such as $P' \rightarrow c$, $P'' \rightarrow c$, and $P^c \rightarrow c$, are redundant since they have the same support and interestingness as rule $P \rightarrow c$ but are more specific. A rule set is non-redundant if it includes all rules except redundant rules.

Theorem 2 *The relationship with the non-redundant rule set*

An optimal rule set is a subset of a non-redundant rule set.

A proof is provided in the appendix.

We have another special case for Theorem 1.

Corollary 2 *Termination property*

If $\text{supp}(P \neg c) = 0$, then all more specific rules of the rule $P \rightarrow c$ do not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Klogen, conviction, p -s (or leverage), Laplace, cosine, certainty factor or Jaccard.

A proof is provided in the appendix.

The practical implication of the above corollary is that: once $\text{supp}(P \neg c) = 0$ is observed, it is not necessary to search for more specific rules of $P \rightarrow c$. Those more specific rules will not be in an optimal rule set. Rule $P \rightarrow c$ is kept since it may be in an optimal rule set.

Let us look at why it is called the termination property. Assume that the interestingness metric is confidence. If $\text{supp}(P \neg c) = 0$ then rule $\text{conf}(P \rightarrow c) = 100\%$. None of its more specific rules improves this confidence, and we stop going any further. The

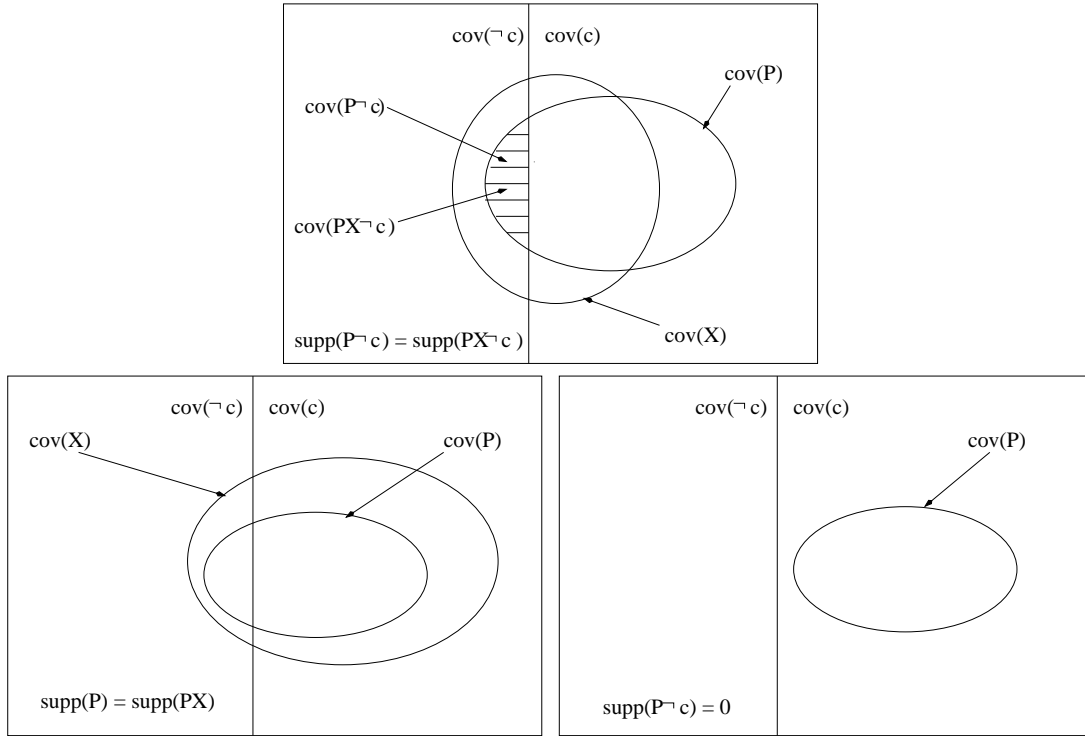


Figure 1: Depictions of conditions of Theorem 1, and Corollaries 1 and 2. In each diagram, the left hand side rectangle stands for the set of records not containing the class c , and the right hand side rectangle stands for the set of records containing the class c . Ellipses denote the cover sets of patterns P and X

same is true for other interestingness metrics.

An illustrative comparison of conditions of Theorem 1, and Corollaries 1 and 2 is given in Figure 1.

We use the following example to show how the Theorem 1 and its corollaries work.

Example 3 We assume z is a class and do not assume the minimum support requirement in this example.

	b	c		z
a	b	c	d	z
a		c		z
a	b		d	z
	b		d	$\neg z$
a	b	c		$\neg z$

Usually, we have to consider all 15 candidate rules: $a \rightarrow z$, $b \rightarrow z$, $c \rightarrow z$, $d \rightarrow z$, $ab \rightarrow z$, $ac \rightarrow z$, $ad \rightarrow z$, $bc \rightarrow z$, $bd \rightarrow z$, $cd \rightarrow z$, $abc \rightarrow z$, $abd \rightarrow z$, $acd \rightarrow z$, $bcd \rightarrow z$, and $abcd \rightarrow z$.

Since $\text{supp}(ab\neg z) = \text{supp}(a\neg z)$, according to Theorem 1 rule $ab \rightarrow z$ and all its more specific rules, $abc \rightarrow z$, $abd \rightarrow z$ and $abcd \rightarrow z$ will not be in an optimal rule set. Similarly, both rule $ac \rightarrow z$ and rule $bc \rightarrow z$ and their more specific rules will not be in the optimal rule set.

Since $\text{supp}(ad\neg z) = 0$, all more specific rules of $ad \rightarrow z$, e.g. $abd \rightarrow z$, $acd \rightarrow z$ and $abcd \rightarrow z$ will not be in the optimal rule set. Similarly, all more specific rules of $cd \rightarrow z$ will not be in the optimal rule set.

Since $\text{supp}(d) = \text{supp}(bd)$, $bd \rightarrow c$ and all its more specific rules will not in the optimal rule set.

Therefore, rules that are possible to be in the optimal rule set are $a \rightarrow z$, $b \rightarrow z$, $c \rightarrow z$, $d \rightarrow z$, $ad \rightarrow z$ and $cd \rightarrow z$. This set of candidate rules is significantly smaller than the original 15 candidate rules.

4 An optimal rule discovery algorithm

In this section, we first show how to use Theorem 1 and its corollaries for forward pruning. Then we discuss a candidate presentation method for easy pruning. Subsections 4.3 and 4.4 present the detailed implementation of forward pruning by Theorem 1 and its corollaries. The complete algorithm is presented in Subsection 4.5. After that, an illustrative example shows how the algorithm works. Finally, we discuss some implementation issues such as counting and data structure.

4.1 Basic ideas and forward pruning

General to specific searching is very common in rule discovery. For example, C4.5rules, CN2 and Apriori employ this method.

When a heuristic search method is employed, we worry less about combinatorial explosion. However, when we conduct optimal search, combinatorial explosion is a big problem.

We first look at how association rule discovery alleviates this problem. An association rule discovery algorithm prunes infrequent patterns forwardly. A pattern is frequent if its support is not less than the minimum support. A pattern is potentially frequent only if all its sub patterns are frequent, and this anti-monotonic property is used to limit the number of patterns to be searched. This is called forward pruning.

Optimal rule discovery makes use of Theorem 1 and Corollaries to forwardly prune rules.

We illustrate how Theorem 1 forwardly prunes rules that are not in an optimal set. Given a pattern $abcd$, and assume the target is fixed to z . We usually have to examine candidate rules $a \rightarrow z, b \rightarrow z, \dots$ for 1-patterns, $ab \rightarrow z, ac \rightarrow z, \dots$ for 2-patterns, $abc \rightarrow z, abd \rightarrow z, \dots$ for 3-patterns, and $abcd \rightarrow z$ for 4-pattern. If we know $\text{supp}(a \rightarrow z) = \text{supp}(ab \rightarrow z)$, then the theorem empower us to skip examining candidates, $ab \rightarrow z, abc \rightarrow z, abd \rightarrow z$, and $abcd \rightarrow z$ because they are not in the optimal rule set anyway. Corollary 2 works in the similar way.

We show how Corollary 1 forwardly prunes rules that are not in an optimal rule set by using the above example. When the targets are not fixed to z but include x and y too, the number of rule candidates is tripled. We list those including pattern ab in their antecedents as follows, $ab \rightarrow x, ab \rightarrow y, ab \rightarrow z, abc \rightarrow x, abc \rightarrow y, abc \rightarrow z, abd \rightarrow x, abd \rightarrow y, abd \rightarrow z, abcd \rightarrow x, abcd \rightarrow y$, and $abcd \rightarrow z$. If we know that $\text{supp}(a) = \text{supp}(ab)$ holds, then all above listed candidates are ignored according to

Corollary 1. Corollary 1 prunes more candidates than Theorem 1, but its requirement is stricter.

4.2 Candidate representation

To facilitate the implementation of forward pruning by Theorem 1 and its corollaries, we define a rule candidate as a pair of (pattern, target-set), denoted by (P, C) . P is a pattern and C is a set of classes. In a relational data set, we have $P \cap C = \emptyset$. For example, let $P = abc$ and $C = xyz$, and then (P, C) stands for three candidate rules, $abc \rightarrow x$, $abc \rightarrow y$ and $abc \rightarrow z$. To remove a candidate rule, we simply remove a class from the target-set. For example, candidate $\{abc, yz\}$ stands for only two candidate rules, namely $abc \rightarrow y$ and $abc \rightarrow z$. When the target-set is empty, the candidate stands for no rules.

The removal of classes from target-set C is determined by Theorem 1 and its corollaries. For example, if we have $\text{supp}(ab\neg x) = \text{supp}(abc\neg x)$, then according to the Theorem, rule $abc \rightarrow x$ and all its more specific rules will not occur in the optimal rule set. x should be removed from the candidate set and the candidate becomes (abc, yz) . If we know $\text{supp}(abc\neg y) = 0$, y should be removed from the target-set according to Corollary 2 and the candidate becomes (abc, z) . Consider another candidate (bcd, xyz) . If we have $\text{supp}(bcd) = \text{supp}(cd)$, then according to Corollary 1 the target-set of the candidate should be emptied and the candidate becomes (bcd, \emptyset) .

Candidate (P, \emptyset) should be removed since no rules will be generated from it and its super candidates. (P', C') is called a super candidates of (P, C) if $P' \supset P$ holds.

The existence of a candidate relies on two conditions: (1) pattern P is frequent, and (2) target-set C is not empty.

4.3 Candidate generator

For the easy comparison, we present Candidate-gen for the optimal rule set discovery in the similar way as does the Apriori-gen. We call a candidate l -candidate if its pattern is a l -pattern. An l -candidate set includes all l -candidates.

Function Candidate-gen

- ```
// Combining
```
- 1) for each pair of candidates  $(P_{l-1}s, C_s)$  and  $(P_{l-1}t, C_t)$  in  $l$ -candidate set
  - 2) insert candidate  $(P_{l+1}, C)$  in the  $(l + 1)$ -candidate set  
where  $P_{l+1} = P_{l-1}st$  and  $C = C_s \cap C_t$
- ```
// Pruning
```
- 3) for all $P_l \subset P_{l+1}$
 - 4) if candidate (P_l, C_l) does not exist
 - 5) then remove candidate (P_{l+1}, C) and return
 - 6) else $C = C \cap C_l$
 - 7) If the target-set of (P_{l+1}, C) is empty
 - 8) then remove the candidate

We first explain lines 1 and 2. Suppose that we have two candidates (abc, xy) and (abd, y) . The new candidate is $(abcd, y)$. The intersection of target-sets here and in line 6 is to ensure that removed classes from the target-set of a candidate never appear in the target-set of its super candidates. The correctness is guaranteed by Theorem 1 and its corollaries since any class removal in the target-set is determined by them.

Second we explain lines 3 to 8. Suppose that we have new candidate $(abcd, y)$. It is the combination of (abc, y) and (abd, yz) . We need to check if candidates identified by patterns $\{acd\}$ and $\{bcd\}$ exist. Suppose that they do exist and are (acd, y) and (bcd, xz) . After considering candidate (acd, y) by line 7, the new candidate remains unchanged. After considering candidate (bcd, xz) by line 7, the target-set of the new

candidate becomes empty because of $C = \{y\} \cap \{x, z\} = \emptyset$. The new candidate is then pruned.

4.4 More pruning

We have a pruning process in candidate generation, and will have another pruning process after counting the support of candidates. This is a key to make use of Theorem 1 and its corollaries for pruning. In the following algorithm, σ is the minimum support.

Function Prune($l + 1$)

// $l + 1$ is the new level where candidates are counted.

- 1) for each candidate (P, C) in $(l + 1)$ -candidate set
- 2) for each $c \in C$
 - // test the frequency individually
- 3) if $\text{supp}(Pc) / \text{supp}(c) \leq \sigma$ then remove c from C
 - // test the satisfaction of Corollary 2
- 4) else if $\text{supp}(P\neg c) = 0$ then mark c terminated
- 5) if C is empty then remove candidate (P, C) and return
- 6) for each l level subset $P' \subset P$
 - // test the satisfaction of Corollary 1
- 7) if $\text{supp}(P) = \text{supp}(P')$ then empty C
 - // test the satisfaction of Theorem 1
- 8) else if $\text{supp}(P\neg c) = \text{supp}(P'\neg c)$ then remove c from C
- 9) if C is empty then remove candidate (P, C)

We prune candidates from two aspects, the infrequency of the pattern and the emptiness of the target-set.

Here we consider a local support instead of the global support. Because of the possible skewed distributions of classes, a single global support is not suitable for a variant rules targeting different classes. Many applications have adopted local support in spite of using different names, such as coverage in [3]. The local support of rule $P \rightarrow c$ is defined as $\text{supp}(Pc)/\text{supp}(c)$. The local support is the support in the sub data set containing c . It is also called the recall of rule $P \rightarrow c$. We prefer local support since it observes the anti-monotonic property of the support. When we make use of local support, infrequent candidate rules are removed one by one as in line 3.

We are aware of two variants of the forward pruning by Theorem 1 and Corollary 2. One is that rule $PX \rightarrow c$ and all its more specific rules are not in an optimal rule set as in Theorem 1. In this case, we just remove c from the target set. Another is that all more specific rules of rule $P \rightarrow c$ are not in an optimal rule set except rule $P \rightarrow c$ as in Corollary 2. In this case, we cannot remove c since otherwise we may lose rule $P \rightarrow c$. We design a special statue for this situation, namely termination of target c .

Definition 4 *Target $c \in C$ is terminated in candidate (P, C) if $\text{supp}(Pc) = 0$.*

Terminated c is removed after the rule forming in line 9 of the ORD algorithm.

Line 4 of the Prune function is a direct application of Corollary 2. Target c is marked as terminated and will be removed after a rule is formed. Once c is removed from C , all more specific rules of $P \rightarrow c$ will be removed in the following rule candidate generation.

Line 7 of the Prune function is a direct result from Corollary 1. All classes in C are removed and as a results candidate (P, C) is removed. No super candidates of (P, C) will be formed in the following rule candidate generation.

Line 8 of the Prune function is a direct utilisation of Theorem 1. Target c is removed, and therefore rule $P \rightarrow c$ and all its more specific rules are removed in the following candidate generation.

All candidates with an empty target set are removed in lines 5 and 9 to ensure their super candidates are pruned in the following candidate generation.

4.5 ORD algorithm

Now we are able to present our algorithm for the optimal rule discovery, in abbreviation ORD. In this algorithm, any interestingness metric discussed in Section 3 can be used as a rule selection criterion, and the output rule set is an optimal rule set defined by the interestingness metric. It differs from association rule discovery in that it does not form rules from the set of all frequent patterns but generates rules using partial frequent patterns.

Algorithm 1 *ORD: Optimal Rule Discovery*

Input: a data set D , the minimum local support σ and the minimum interestingness θ by a metric.

Output: an optimal rule set R defined by an interestingness metric

- 1) Set $R = \emptyset$
- 2) Count support of 1-patterns by arrays
- 3) build 1-candidate sets
- 4) Form and add rules to R
- 5) Generate 2-candidate set
- 6) While new candidate set is not empty
- 7) Count support of patterns for new candidates
- 8) Prune candidates in the new candidate set
- 9) Form rules and add optimal rules to R
- 10) Generate next level candidate set
- 11) Return the rule set R

The above algorithm is self explanatory and its two main functions have been discussed in the previous subsections.

The ORD algorithm is efficient since it does not generate all frequent patterns. It only makes use of a small subset of frequent patterns as shown in experiments. Note that line 8 in Function Candidate-gen and line 5 and 9 in Function Prune, all super candidates of a candidate with the empty target-set are removed on top of those of a candidate with the infrequent pattern.

4.6 An illustrative example

We provide an example to show how the ORD algorithm works in this section.

Example 4 *In the following data set, y and z are classes. We do not assume the minimum support constraint. All candidates generated by the ORD algorithm are listed in Figure 2.*

b	c	d	y
a	c	y	
a	b	c	y
a	b	c	d
a			y
	b	d	z
a	c		z
	b	c	d
a			z
a	b	c	z

The first level candidates are used to prune the second level candidates. z is removed in candidate (ad, yz) by line 3 in Function Prune due to $\text{supp}(adz) = 0$. y in candidate (ad, yz) is terminated by line 4 in Function Prune because of $\text{supp}(ad\bar{y}) =$

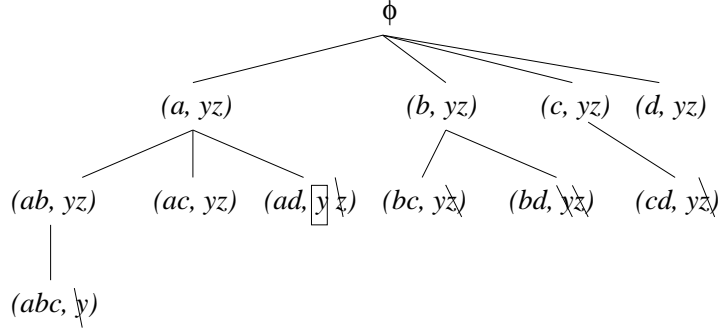


Figure 2: All candidates searched in Example 4. A class crossed is removed and a class boxed is terminated

0. Both y and z are removed in candidate (bd, yz) by line 7 in Function Prune since $\text{supp}(bd) = \text{supp}(d)$ holds. z in candidate (bc, yz) and (cd, yz) is removed by line 8 in Function Prune because both $\text{supp}(bc\neg z) = \text{supp}(b\neg z)$ and $\text{supp}(cd\neg z) = \text{supp}(d\neg z)$ hold.

After rules have been formed, candidates (ad, \emptyset) and (bd, \emptyset) are removed. As a result, all their super candidates, such as, (abd, \emptyset) , (acd, \emptyset) and (bcd, \emptyset) , will not be generated according to lines 4 and 5 in Function Candidate-gen.

Candidate (abc, yz) is generated by combining candidates (ab, yz) and (ac, yz) according to lines 1 and 2 in Function Candidate-gen. z is then pruned by line 7 in Function Candidate Generator using its sub candidate (bc, y) . y in candidate (abc, y) is removed by line 8 in Function Prune due to $\text{supp}(abc\neg y) = \text{supp}(ab\neg y)$. Subsequently, candidate (abc, \emptyset) is removed.

In rule forming procedure, rules are formed by a user specified interestingness metric. No matter what a metric discussed in 3 is used, the set of candidates is identical. Only the output rule set differs.

5 Support pruning, closure pruning and optimality pruning

In this section we discuss support pruning, closure pruning, and optimality pruning and then characterise relationships among them. This clarifies the efficiency improvement of optimal rule discovery over association rule discovery and non-redundant rule discovery.

First look at the support pruning of the following data set.

<i>b c d e</i>
<i>a b d</i>
<i>a b c</i>
<i>a b c d</i>
<i>a c</i>
Data set A

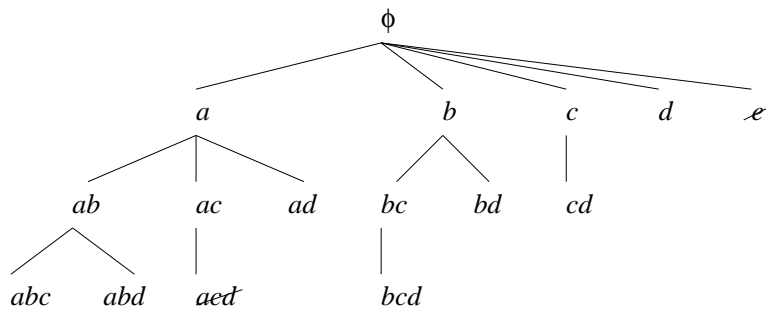


Figure 3: Support pruning for mining frequent patterns on data set A. Patterns crossed are infrequent

Figure 3 shows the support pruning by using the minimum support of 0.2. We see that the removal of e in Level 1 equals the removal of 15 patterns in the subsequent levels, such as $ae, be, \dots, abe, ace, \dots, abce, abde, \dots$ and $abcde$.

Support pruning works effectively when the underlying data set is sparse or the minimum support is high. However, it does not work well on dense data sets or when

the minimum support is low.

Look at the closure pruning by the following data set.

<i>a</i>	<i>b</i>	<i>d</i>	
<i>a</i>	<i>c</i>	<i>e</i>	
	<i>b</i>	<i>e</i>	
		<i>d</i>	<i>e</i>
Data set B			

Figure 4. shows the closure pruning. There is no minimum support requirement.

As discussed in Section 3, Corollary 1 summarises the closure pruning. The pattern in box are terminated because the support of its super patterns equals that of itself. If the final goal is to find minimal generators, all candidates are left as they are in Figure 4. If the the final goal is to find closed patterns, the number of candidates remains unchanged, but patterns in some candidates are extended. For example, ac is terminated due to $\text{supp}(ac) = \text{supp}(c)$, and as a result, all occurrences of c will be replaced by ac . Similarly, all occurrences of c is further replaced by ce because of the termination of ce by $\text{supp}(ce) = \text{supp}(c)$. Pattern ace is the only closed pattern left out in Figure 4, and other minimal generators are closed patterns too. Closed patterns can be discovered in the same search tree finding minimal generators. Therefore, both closed pattern mining and minimal generator mining search for the same number of candidates and make use of the closure pruning strategy.

The closure pruning works effectively when the underlying data set is dense or the minimum support is low. We use an example to elaborate the first point. Assume that a dense data set contains five identical records $\{a, b, c, d, e\}$. Closed pattern mining will stop at level 2 after searching for 15 candidates. In contrast, frequent pattern mining will continue all the way to level 5 and search for $2^5 - 1$ candidates. We present the following justification for the second point: $\text{supp}(PX) = \text{supp}(P)$ means $\text{cov}(P) \subseteq$

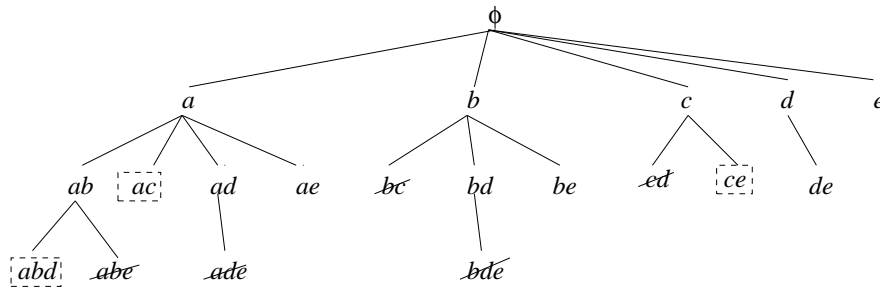


Figure 4: Closure pruning for mining minimal generators on data set B. Patterns crossed are non-existing and patterns boxed are terminated

$cov(X)$. When $cov(X)$ remains unchanged, pattern P with a smaller $cov(P)$ is more probable to satisfy $cov(P) \subseteq cov(X)$ than with a larger $cov(P)$.

The above two pruning strategies are complementary. How does the ORD algorithm use them in an effective way?

Let us look at the following data set concatenating the above two data sets. z is the target for rules. Figure 5 shows the optimality pruning with the minimum local support of 0.2 in sub data set containing z .

b	c	d	e	z
a	b	d		z
a	b	c		z
a	b	c	d	z
a	c			z
a	b	d		$\neg z$
a	c	e		$\neg z$
b		e		$\neg z$
		d	e	$\neg z$

Data set (A + B)

The candidate set searched by optimal rule discovery is the intersection of the candidate set (frequent patterns) for association rule discovery in z sub data set and the

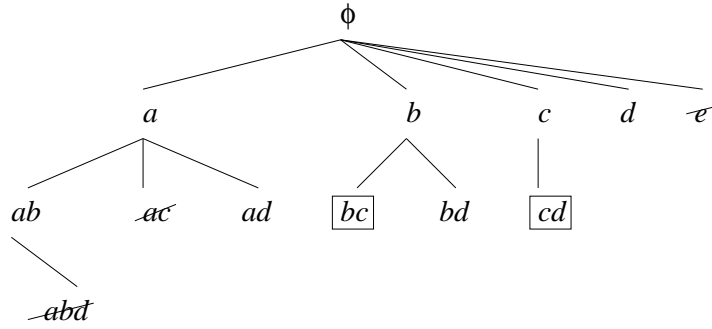


Figure 5: Optimality pruning for optimal rule discovery targeting z on data set $(A+B)$. Patterns crossed are removed and patterns boxed are terminated

candidate set (minimal generators) for non-redundant rule discovery in $\neg z$ sub data set. The crucial point is that both have to perform simultaneously. Both association rule discovery and non-redundant rule discovery search for more candidates than optimal rule discovery.

Now, we have an insight understanding of optimality pruning stated in Theorem 1. It makes use of the closure pruning strategy. Comparing Corollary 1 with Theorem 1, we find that Theorem 1 states Corollary 1 in the sub data set excluding c .

6 Experimental results

In this section, we empirically evaluate the computational complexity of optimal rule discovery in comparison with association rule discovery and non-redundant association rule discovery on four data sets from UCML repository [4] described in Table 1. The efficiency of optimal rule discovery is its effective optimality pruning. We show that the optimality pruning significantly reduces candidates for searching.

The efficiency of an algorithm depends significantly on the data structure and implementation. For example, association rule discovery has various implementations. All are based on support pruning strategy, but their execution time varies. Theoretically, their computational complexities are the same since they all search for all fre-

Name	#Records	#attribute	#classes
Anneal	898	38	5
Hypothyroid	3163	25	2
Mushroom	8124	22	2
Sick	2800	29	2

Table 1: A brief description of data sets

quent patterns. Their efficiencies vary since they employ different data structures and counting schemes. There are a number of implementations for association rule discovery, and we are unable to compare with them individually by execution time.

However, the computational complexity improvement is fundamental and the implementation only accelerates the improvement. An empirical estimation of the complexity for a rule discovery algorithm is the number of candidates it searches. In this paper, we compare the searched candidates for association rule discovery, for non-redundant association rule discovery, and for optimal rule discovery. An association discovery algorithm searches for all frequent patterns, and a non-redundant rule discovery algorithm searches for all frequent minimal generators (equivalently all frequent closed patterns). We compare the number of candidates for the ORD algorithm with the number of frequent patterns and the number of frequent minimal generators. This comparison is independent of the implementation.

In this experiment, we employ the local support as defined in Subsection 4.4, which is a ratio for an individual class. A pattern is frequent if it is frequent in at least one class. All frequent patterns are stored in the prefix tree.

The experiment was conducted on a 1 GHz CPU computer with 2G memory running Linux. We search for rules containing up to eight attribute-value pairs. We do not specify the type of optimal rule set since the same candidate set generates an optimal rule set defined by any interestingness metric discussed in Section 3.

Figure 6 shows the searched candidates by the ORD algorithm in comparison with

the number of frequent patterns and the number of frequent minimal generators. The set of ORD candidates is a very small subset of frequent patterns, and a subset of frequent minimal generators. This trend is more evident when the minimum support is low. This shows that optimal rule discovery has significant less computational complexity than association rule discovery, and less computational complexity than non-redundant association rule discovery.

In comparison with optimal rule discovery and non-redundant rule discovery, association rule discovery is very inefficient in data sets like ones used in this experiment. The efficiency of association rule discovery deteriorates dramatically when the minimum support is low. Optimal rule discovery is more efficient than non-redundant association rule discovery. Though differences between candidate numbers of non-redundant association rule discovery and optimal rule discovery are squashed in Figure 6 by the large number of frequent patterns, the discrepancies are still clear in data sets Mushrooms and Sick.

7 Conclusions

In this paper, we discussed a family of optimal rule sets, the properties for their efficient discovery and their relationships with the non-redundant rule sets. The family of optimal rule sets support a simple anti-monotonic property and an optimal rule set is a subset of a non-redundant rule set. We presented the ORD algorithm for mining optimal rule sets, and evaluated its computational complexity on some data sets in comparison with the association rule discovery and non-redundant association rule discovery. The computational complexity of optimal rule discovery is significantly lower than that of association rule discovery and lower than that of non-redundant association rule discovery. We discussed the relationship of optimal pruning with support pruning and closure pruning, and concluded that optimality pruning makes use of both

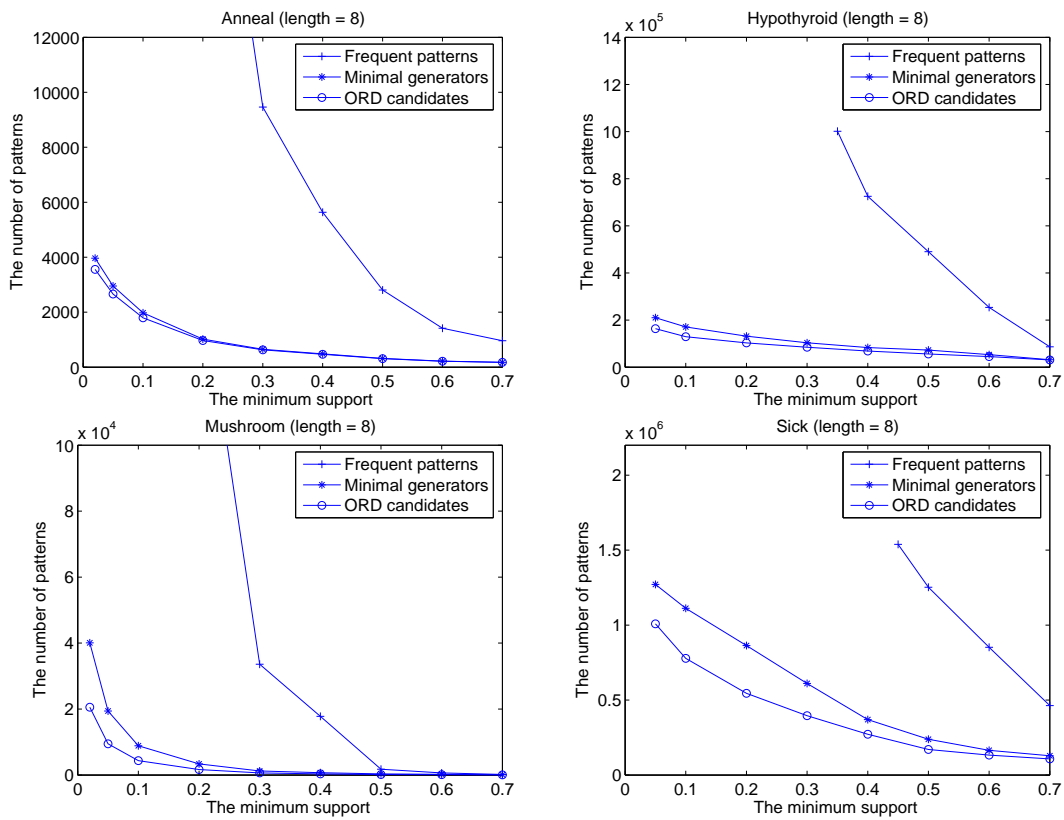


Figure 6: The number of candidates for optimal rule discovery versus the number of frequent patterns for association rule discovery and the number of frequent minimal generators for non-redundant association rule discovery. The ORD searches a small subset of frequent patterns, and a subset of minimal generators.

support and closure pruning strategies simultaneously on two disjointed data sets.

Optimal rule discovery is efficient and works well with the low or no minimum support constraint. It generates optimal rule sets for a number of interestingness metrics. Therefore, it is a great alternative for association rule discovery.

Acknowledgements

I thank Prof Tony Roberts, Prof Ada Fu, and Dr Xiaodi Huang for their constructive comments. This project is supported by ARC (Australia Research Council) grant DP0559090.

References

- [1] R. Agrawal, T. Imielinski, and A. Swami. Mining associations between sets of items in massive databases. In *Proc. of the ACM SIGMOD Int'l Conference on Management of Data*, pages 207–216, 1993.
- [2] R. Bayardo and R. Agrawal. Mining the most interesting rules. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 145–154, N.Y., 1999. ACM Press.
- [3] R. Bayardo, R. Agrawal, and D. Gunopulos. Constraint-based rule mining in large, dense database. *Data Mining and Knowledge Discovery Journal*, 4(2/3):217–240, 2000.
- [4] E. K. C. Blake and C. J. Merz. UCI repository of machine learning databases, <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998.
- [5] S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. In *Proceedings, ACM SIGMOD International Conference on Management of Data: SIGMOD 1997: May 13–15, 1997, Tucson, Arizona, USA*, volume 26(2), pages 255–264, NY, USA, 1997. ACM Press.
- [6] P. Clark and R. Boswell. Rule induction with CN2: Some recent improvements. In *Machine Learning - EWSL-91*, pages 151–163, 1991.

- [7] W. W. Cohen. Fast effective rule induction. In *Proceedings of the Twelfth International Conference on Machine Learning (ICML)*, pages 115–123, 1995.
- [8] V. Dhar and A. Tuzhilin. Abstract-driven pattern discovery in databases. *IEEE Transactions on Knowledge and Data Engineering*, 5(6), 1993.
- [9] T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Data mining using two-dimensional optimized association rules: scheme, algorithms, and visualization. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data, Montreal, Quebec, Canada, June 4–6, 1996*, pages 13–23, New York, 1996. ACM Press.
- [10] H. Hu and J. Li. Using association rules to make rule-based classifiers robust. In *Proceedings of Sixteenth Australasian Database Conference (ADC)*, page 47–52, Newcastle Australia, 2005. ACS Society.
- [11] J. Li, H. Shen, and R. Topor. Mining the optimal class association rule set. *Knowledge-Based System*, 15(7):399–405, 2002.
- [12] B. Liu, W. Hsu, and Y. Ma. Integrating classification and association rule mining. In *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining (KDD-98)*, pages 27–31, 1998.
- [13] E. Omiecinski. Alternative interest measures for mining associations in databases. *IEEE Transactions on Knowledge and Data Engineering*, 15(1), 2003.
- [14] G. Piatetsky-Shapiro. Discovery, analysis and presentation of strong rules. In G. Piatetsky-Shapiro, editor, *Knowledge Discovery in Databases*, pages 229–248. AAAI Press / The MIT Press, Menlo Park, California, 1991.
- [15] J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, CA, 1993.
- [16] P. Tan, V. Kumar, and J. Srivastava. Selecting the right objective measure for association analysis. *Information Systems*, 29(4):293 – 313, 2004.
- [17] G. I. Webb. OPUS: An efficient admissible algorithm for unordered search. In *Journal of Artificial Intelligence Research*, volume 3, pages 431–465, 1995.
- [18] G. I. Webb and S. Zhang. K-optimal rule discovery. In *Data Mining and Knowledge Discovery Journal*, volume 10(1), pages 39 – 79, 2005.
- [19] M. J. Zaki. Mining non-redundant association rules. *Data Mining and Knowledge Discovery Journal*, 9:223–248, 2004.
- [20] M. J. Zaki and C. J. Hsiao. Charm: An efficient algorithm for closed association rule mining. In *Proceedings of SIAM International Conference on Data Mining*, 2002.

Appendix

In this Appendix, we provide proofs for Theorems 1 and 2, and Corollaries 1 and 2.

Theorem 1 Anti-monotonic property

if $\text{supp}(PX \neg c) = \text{supp}(P \neg c)$ then rule $PX \rightarrow c$ and all its more specific rules will not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Klossgen, conviction, p-s (or leverage), Laplace, cosine, certainty factor or Jaccard.

Proof In the proof, we show $\text{Interestingness}(PQX \rightarrow c) \leq \text{Interestingness}(PQ \rightarrow c)$. Therefore, rule $(PX \rightarrow c)$ (when $Q = \emptyset$) and all its more specific rules, for example $PQX \rightarrow c$ (when $Q \neq \emptyset$), will not occur in the optimal rule set.

The only case for the condition $\text{supp}(PX \neg c) = \text{supp}(P \neg c)$ holding is that $\text{cov}(P \neg c) \subseteq \text{cov}(X \neg c)$. We then deduce that $\text{cov}(PQ \neg c) \subseteq \text{cov}(QX \neg c)$ for any Q . consequently, $\text{supp}(PQX \neg c) = \text{supp}(PQ \neg c)$ holds for any Q .

For the confidence case, consider $f(y) = y/(y + \alpha)$ monotonically increases with y when constant $\alpha > 0$ and $\text{supp}(PQ) \geq \text{supp}(PQX) > 0$:

$$\begin{aligned}
 \text{conf}(PQ \rightarrow c) &= \frac{\text{supp}(PQc)}{\text{supp}(PQ)} \\
 &= \frac{\text{supp}(PQc)}{\text{supp}(PQc) + \text{supp}(PQ \neg c)} \\
 &= \frac{\text{supp}(PQc)}{\text{supp}(PQc) + \text{supp}(PQX \neg c)} \\
 &\geq \frac{\text{supp}(PQXc)}{\text{supp}(PQXc) + \text{supp}(PQX \neg c)} \\
 &= \text{conf}(PQX \rightarrow c).
 \end{aligned}$$

We then prove the odds ratio case. Odds ratio is a classic statistical metric to

measure the association between events. Consider $f(y) = y/(\alpha - y)$ monotonically increases with y when constant $\alpha > 0$ and $\text{supp}(PQ) \geq \text{supp}(PQX) > 0$:

$$\begin{aligned}
\text{or}(PQ \rightarrow c) &= \frac{\text{supp}(PQc) \text{supp}(\neg(PQ)\neg c)}{\text{supp}(\neg(PQ)c) \text{supp}(PQ\neg c)} \\
&= \frac{\text{supp}(PQc)(\text{supp}(\neg c) - \text{supp}(PQ\neg c))}{(\text{supp}(c) - \text{supp}(PQc)) \text{supp}(PQ\neg c)} \\
&= \frac{\text{supp}(PQc)(\text{supp}(\neg c) - \text{supp}(PQX\neg c))}{(\text{supp}(c) - \text{supp}(PQc)) \text{supp}(PQX\neg c)} \\
&= \frac{\text{supp}(PQc) \text{supp}(\neg(PQX)\neg c)}{(\text{supp}(c) - \text{supp}(PQc)) \text{supp}(PQX\neg c)} \\
&\geq \frac{\text{supp}(PQXc) \text{supp}(\neg(PQX)\neg c)}{(\text{supp}(c) - \text{supp}(PQXc)) \text{supp}(PQX\neg c)} \\
&= \frac{\text{supp}(PQXc) \text{supp}(\neg(PQX)\neg c)}{\text{supp}(\neg(PQX)c) \text{supp}(PQX\neg c)} \\
&= \text{or}(PQX \rightarrow c).
\end{aligned}$$

Lift also known as interest [5] or strength [8], is a widely used metric for ranking the interestingness of association rules. It has been used in IBM Intelligent Miner. We make use of the previous results, $\text{conf}(PQ \rightarrow c) \geq \text{conf}(PQX \rightarrow c)$, in the following proof:

$$\begin{aligned}
\text{lift}(PQ \rightarrow c) &= \frac{\text{supp}(PQc)}{\text{supp}(PQ) \text{supp}(c)} \\
&= \frac{\text{conf}(PQ \rightarrow c)}{\text{supp}(c)} \\
&\geq \frac{\text{conf}(PQX \rightarrow c)}{\text{supp}(c)} \\
&= \text{lift}(PQX \rightarrow c).
\end{aligned}$$

Gain [9] is an alternative for confidence. Fraction θ is a constant in interval $(0, 1)$, and only rules obtaining positive gain are interesting. We use $\text{conf}(PQ \rightarrow c) \geq$

$\text{conf}(PQX \rightarrow c) > \theta$ in the following proof:

$$\begin{aligned}
\text{gain}(PQ \rightarrow c) &= \text{supp}(PQc) - \theta \text{supp}(PQ) \\
&= (\text{conf}(PQ \rightarrow c) - \theta) \text{supp}(PQ) \\
&\geq (\text{conf}(PQX \rightarrow c) - \theta) \text{supp}(PQX) \\
&= \text{gain}(PQX \rightarrow c).
\end{aligned}$$

The proofs for metrics added-value, $\text{addedvalue}(P \rightarrow c) = \text{conf}(P \rightarrow c) - \text{supp}(c)$, and Klogen, $\text{Klogen}(P \rightarrow c) = \sqrt{\text{supp}(Pc)}(\text{conf}(P \rightarrow c) - \text{supp}(c))$, are very straightforward and hence we omit them here.

Conviction [5] is used to measure deviations from the independence by considering outside negation:

$$\begin{aligned}
\text{conviction}(PQ \rightarrow c) &= \frac{\text{supp}(PQ) \text{supp}(\neg c)}{\text{supp}(PQ\neg c)} \\
&= \frac{\text{supp}(PQ)(1 - \text{supp}(c))}{\text{supp}(PQ) - \text{supp}(PQc)} \\
&= \frac{1 - \text{supp}(c)}{1 - \text{conf}(PQ \rightarrow c)} \\
&\geq \frac{1 - \text{supp}(c)}{1 - \text{conf}(PQX \rightarrow c)} \\
&= \text{conviction}(PQX \rightarrow c).
\end{aligned}$$

P-s metric (or leverage), $\text{ps}(P \rightarrow c) = \text{supp}(Pc) - \text{supp}(P) \text{supp}(c)$, is a classic interestingness metric for rules proposed by Piatesky-Shaprio [14]. The proof for it is very similar to that of gain and hence we omit it.

Laplace [6, 17] accuracy is a metric for classification rules. $|D|$ is the number of transactions in D and k is the number of classes. In classification problems, $k \geq 2$ and usually $\text{conf}(PQX \rightarrow c) \geq 0.5$. Therefore, $k \cdot \text{conf}(PQX \rightarrow c) \geq 1$ holds. Function

$f(y) = (\alpha|D| + y)/(|D| + ky)$ monotonically decrease with y when $k \cdot \alpha > 1$ and $1/\text{supp}(PQX) \geq 1/\text{supp}(PQ)$.

$$\begin{aligned}
\text{Laplace}(PQ \rightarrow c) &= \frac{\text{supp}(PQc)|D| + 1}{\text{supp}(PQ)|D| + k} \\
&= \frac{\text{conf}(PQ \rightarrow c)|D| + 1/\text{supp}(PQ)}{|D| + k/\text{supp}(PQ)} \\
&\geq \frac{\text{conf}(PQX \rightarrow c)|D| + 1/\text{supp}(PQ)}{|D| + k/\text{supp}(PQ)} \\
&\geq \frac{\text{conf}(PQX \rightarrow c)|D| + 1/\text{supp}(PQX)}{|D| + k/\text{supp}(PQX)} \\
&= \text{Laplace}(PQX \rightarrow c).
\end{aligned}$$

The proofs for the following two metrics are straightforward and hence we omit them. $\text{Cosine}(P \rightarrow c) = \text{supp}(Pc)/(\sqrt{\text{supp}(P)\text{supp}(c)})$ and $\text{Certaintyfactor}(P \rightarrow c) = (\text{conf}(P \rightarrow c) - \text{supp}(c))/(1 - \text{supp}(c))$.

Finally, we prove the metric of Jaccard.

$$\begin{aligned}
\text{Jaccard}(PQ \rightarrow c) &= \frac{\text{supp}(PQc)}{\text{supp}(PQ) + \text{supp}(c) - \text{supp}(PQc)} \\
&= \frac{\text{conf}(PQ \rightarrow c)}{1 + \text{supp}(c)/\text{supp}(PQ) - \text{conf}(PQ \rightarrow c)} \\
&\geq \frac{\text{conf}(PQ \rightarrow c)}{1 + \text{supp}(c)/\text{supp}(PQX) - \text{conf}(PQ \rightarrow c)} \\
&\geq \frac{\text{conf}(PQX \rightarrow c)}{1 + \text{supp}(c)/\text{supp}(PQX) - \text{conf}(PQX \rightarrow c)} \\
&= \text{Jaccard}(PQX \rightarrow c).
\end{aligned}$$

The theorem has been proved. \square

Corollary 1 Closure property

If $\text{supp}(P) = \text{supp}(PX)$, then rule $PX \rightarrow c$ for any c and all its more specific rules

do not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Kloggen, conviction, p-s (or leverage), Laplace, cosine, certainty factor or Jaccard.

Proof If $\text{supp}(P) = \text{supp}(PX)$, then $\text{supp}(P \neg c) = \text{supp}(PX \neg c)$ holds for any c . Therefore, this Corollary is proved immediately by Theorem 1. \square

Corollary 2 Termination property

If $\text{supp}(P \neg c) = 0$, then all more specific rules of the rule $P \rightarrow c$ do not occur in an optimal rule set defined by confidence, odds ratio, lift (interest or strength), gain, added-value, Kloggen, conviction, p-s (or leverage), Laplace, cosine, certainty factor or Jaccard.

Proof If $\text{supp}(P \neg c) = 0$ then $\text{supp}(PX \neg c) = \text{supp}(P \neg c) = 0$ holds for any X . Therefore, this Corollary is proved immediately by Theorem 1. \square

Theorem 2 The relationship with the non-redundant rule set

An optimal rule set is a subset of a non-redundant rule set.

Proof Suppose that we have $\text{supp}(P) = \text{supp}(PX)$ and there is no $P' \subset P$ such that $\text{supp}(P) = \text{supp}(P')$. The rule $PX \rightarrow c$ for any c is redundant. It will not be in an optimal rule set either according to Corollary 1.

Suppose that $\text{supp}(P) = \text{supp}(PX)$ and there is $P' \subset P$ such that $\text{supp}(P) = \text{supp}(P')$. We always have $\text{supp}(P') = \text{supp}(P'Y)$ for $Y \subseteq (PX \setminus P')$. Rules $P \rightarrow c$ and $PX \rightarrow c$ are redundant. They will not be in an optimal rule set either according to Corollary 1.

Further, many non-redundant rules are pruned by Theorem 1.

Therefore, an optimal rule set is a subset of non-redundant rule set. \square