

Pruning and Summarizing the Discovered Associations

Bing Liu, Wynne Hsu and Yiming Ma

School of Computing
National University of Singapore
Lower Kent Ridge Road, Singapore 119260
{liub, whsu, maym}@comp.nus.edu.sg

Abstract

Association rules are a fundamental class of patterns that exist in data. The key strength of association rule mining is its completeness. It finds all associations in the data that satisfy the user specified minimum support and minimum confidence constraints. This strength, however, comes with a major drawback. It often produces a huge number of associations. This is particularly true for data sets whose attributes are highly correlated. The huge number of associations makes it very difficult, if not impossible, for a human user to analyze in order to identify those interesting/useful ones. In this paper, we propose a novel technique to overcome this problem. The technique first prunes the discovered associations to remove those insignificant associations, and then finds a special subset of the unpruned associations to form a summary of the discovered associations. We call this subset of associations the *direction setting* (DS) rules as they set the directions that are followed by the rest of the associations. Using this summary, the user can focus on the essential aspects (or relationships) of the domain and selectively view the relevant details. The approach is effective because experiment results show that the set of DS rules is typically very small. They can be analyzed manually by a human user. The proposed technique has also been applied successfully to a number of real-life applications.

1. Introduction

Association rules are a class of important regularities in data. Association rule mining is commonly stated as follows [2]: Let $I = \{i_1, \dots, i_n\}$ be a set of *items*, and D be a set of data cases. Each data case consists of a subset of items in I . An *association rule* is an implication of the form $X \rightarrow Y$, where $X \subset I$, $Y \subset I$, and $X \cap Y = \emptyset$. The rule $X \rightarrow Y$ holds in D with confidence c if $c\%$ of data cases in D that support X also support Y . The rule has support s in D if $s\%$ of the data case in D contains $X \cup Y$. The problem of mining association rules is to generate all association rules that have support and confidence greater than the user-specified minimum support and minimum confidence.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD-99 San Diego CA USA

Copyright ACM 1999 1-58113-143-7/99/08...\$5.00

A typical association rule mining algorithm works in two steps. The first step finds all *large* itemsets (a set of items) that meet the minimum support constraint. The second step generates rules from all large itemsets that satisfy the minimum confidence constraint.

The key strength of association rule mining is that it can efficiently discover the complete set of associations that exist in data. These associations provide a complete picture of the underlying regularities in the domain. However, this strength comes with a major drawback. The number of discovered associations can be huge, easily in the thousands or tens of thousands (see Section 8). Clearly, such a large number of associations are very difficult, if not impossible, to be analyzed by a human user. This problem is particularly bad with those data sets whose items are highly correlated.

However, it will not be satisfactory to simply give an arbitrary small subset of the rules to the user if there are indeed a large number of them that exist in the data because this small subset can only give a partial picture of the domain. The question then is: "Can we preserve the full power of association rule mining (i.e., its completeness) without overwhelming the user?" This paper shows that it is possible. A novel technique is proposed to prune those insignificant rules and to find a special subset of association rules that represent the essential underlying relationships in the data. We call this subset of associations the *direction setting* (DS) rules. The DS rules give a summary of the behavior of the discovered associations. They represent the essential relationships or structure (or skeleton) of the domain. The non-DS rules simply give additional details. Using the DS rules as a summary, the user can interactively focus on the key aspects of the domain and selectively view the relevant details (non-DS rules).

In this work, we focus on association rule mining from a relational table, which consists of a set of records described by a number of attributes. An item is an attribute value pair, i.e., (attribute = value) (numeric attributes are discretized). Association rule mining in such data is typically targeted at a specific attribute because the user normally wants to know how other attributes are related to this *target* attribute (which can have many values) [13, 4]. With a target attribute, we can express an itemset as follows (instead of a set of items as in [2]):

$$X \rightarrow y$$

where y is an item (or a value) of the target attribute, and X is a set of items from the rest of the attributes. For

simplicity, we call this itemset a *rule* hereafter, regardless of whether it is significant or not. We also say a rule is *large* if it meets the minimum support.

Note that we do not use minimum confidence in our framework, although we can also use it (see Section 6). Minimum confidence does not reflect the underlying relationship of the domain represented by the data [4]. Instead we use statistical correlation as the basis for finding rules that represent the fundamental relations of the domain.

Figure 1 shows the conceptual flow of the proposed technique¹, which consists of two steps, pruning and summarization. Below, we introduce them briefly.

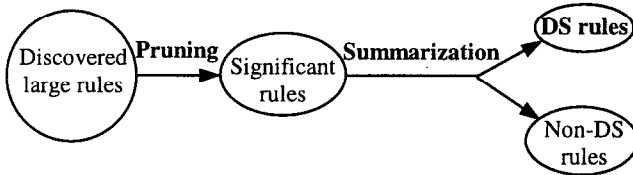


Figure 1. The proposed technique

1.1 Pruning the discovered associations

It is well known that many discovered associations are redundant or minor variations of others. Their existence may simply be due to chance rather than true correlation. Thus, those spurious and insignificant rules should be removed. This is similar to pruning of overfitting rules in classification [21]. Rules that are very specific (with many conditions) tend to overfit the data and have little predictive power. Although association rules are not normally used for prediction, rules that only capture the irregularities and idiosyncrasies of the data have no value and should be removed. An example of such a rule is shown below.

Example 1: We have the following two rules,

- R1: Job = yes \rightarrow Loan = approved
[sup = 60%, conf = 90%]
R2: Job=yes, Credit_history=good \rightarrow Loan= approved
[sup = 40%, conf = 91%]

If we know R1, then R2 is insignificant because it gives little extra information. Its slightly higher confidence is more likely due to chance than to true correlation. It thus should be pruned. R1 is more general and simple. General and simple rules are preferred.

In this work, we measure the significance of a rule using chi-square test (χ^2) for correlation from statistics [17].

1.2 Summarizing the unpruned rules

Pruning can reduce the number of rules substantially. However, the number of rules left can still be very large. This step finds a subset of the rules, called direction setting rules (or DS rules), to summarize the unpruned rules. Essentially, DS rules are significant association rules that set the directions for non-DS rules to follow. The direction of a rule is the type of correlation it has, i.e., *positive correlation* or *negative correlation* or *independence*, which is also computed using χ^2 test. Let us see an example.

Example 2: We have the following discovered rules:

- R1: Job = yes \rightarrow Loan = approved
[sup = 40%, conf = 70%]
R2: Own_house = yes \rightarrow Loan = approved
[sup = 30%, conf = 75%]

χ^2 analysis shows that having a job is positively correlated to the grant of a loan, and owning a house is also positively correlated to obtaining a loan. Then, the following association is not so surprising to us:

- R3: Job = yes, Own_house = yes \rightarrow Loan = approved
[sup = 20%, conf = 90%]

because it intuitively follows R1 and R2. We can use R1 and R2 to provide a summary of the three rules. R1 and R2 are DS rules as they set the direction (positive correlation) that is followed by R3. In real-life data sets, a large number of associations are like R3.

From the example, we see that the DS rules give the essential relationships of the domain. The non-DS rule is not surprising if we already know the DS rules. However, this, by no means, says that non-DS rules are not interesting. Non-DS rules can provide further details about the domain. For example, the non-DS rule above (R3) gives a higher confidence, which may be of interest to the user. Using DS rules to form a summary is analogous to summarization of a text article. From the summary, we know the essence of the article. If we are interested in the details of a particular aspect, the summary can point us to them in the article. In the same way, the DS rules give the essence of the domain and points the user to those related non-DS rules. Non-DS rules are basically combinations of DS rules. In the above Example 2, R3 is a combination of R1 and R2.

Experiment results, including real-life applications, show that although the number of discovered rules can be huge, the number of DS rules is very small (see Section 8). They can be analyzed manually by the human user to obtain the essential relationships in the data. He/she can then focus his/her attention on those interesting aspects of the relationships, and to see the relevant non-DS rules. This process can be easily facilitated by an interactive user interface (Section 7). The proposed technique makes association rule mining effective and practical for data sets whose items are highly correlated. The user can now obtain a complete picture of the domain without being overwhelmed by a huge number of rules.

2. Related Work

The problem of too many rules has been studied by many researchers. [8] proposed an approach to allow the user to specify what he/she wants to see using templates. The system then retrieves those match rules from the set of discovered rules. This method, however, does not prune those insignificant rules and does not provide a summary of the discovered rules.

In subjective interestingness research in data mining, [22, 11, 12, 19] proposed a number of methods for finding unexpected rules. Instead of asking the user to specify what he/she wants to see as in [8], these approaches ask the user

¹ Here, it is presented as a post-processing method. In implementation, it is combined with rule mining.

to specify his/her existing knowledge about the domain. The system then finds those unexpected rules by comparing the user's knowledge with the discovered rules. Again, these methods do not prune and do not attempt to summarize the association rules.

[25] and [18] investigated using item constraints specified by the user in rule mining to generate only the relevant rules. Essentially, the constraints restrict the items or combination of items allowed to participate in mined rules. This approach also does not prune those insignificant rules and does not summarize the unpruned rules.

[26] introduces the concept of association rule cover for pruning association rules. A cover is basically a subset of the discovered associations that can cover the database. The number of rules in a cover can be quite small. A greedy algorithm is proposed to find a good cover and the remaining rules are pruned. The problem with this method is that the advantage of association rules, its completeness, is lost. Clearly, a better approach is to summarize the discovered rules. From this summary, the user can obtain an overall picture of the domain.

[4] proposed a rule pruning technique using *minimum improvement*, which is the difference between the confidence of a rule and the confidence of any proper sub-rule with the same consequent. Those rules that do not meet this minimum improvement in confidence are pruned. [24] also proposed a related technique. Our pruning method is similar. However, we use chi-square test as the basis for pruning (minimum improvement can be easily incorporated in our framework). We will see in Section 8 that even after pruning the number of rules left can still be very large. Summarization is thus important. The methods in [4, 24] do not perform summarization.

There are also a number of other methods in classification research for rule pruning such as pessimistic error rate [21] and minimum description length based pruning [15]. In our work, we choose the widely used chi-square test statistics for rule pruning.

[1] introduces a technique to remove two types of redundant rules, i.e., simple and strict redundancy. Essentially, a rule is redundant with respect to another rule if the support and confidence of the redundant rule are always at least as large as the support and confidence of the latter. Simple redundancy tries to remove those rules that are derived from the same itemset. For example, $AB \Rightarrow C$ is redundant with respect to the rule $A \Rightarrow BC$. This, however, does not happen in our situation because in our case one itemset represents only one rule. Strict redundancy applies to two itemsets and one is a subset of the other, e.g., $X \Rightarrow Y$ is redundant with respect to $X \Rightarrow YZ$. This situation also does not apply in our situation because we focus on association rules that use only one fixed attribute on the right hand side. Also our proposed technique does not use minimum confidence for rule generation (see the problems with minimum confidence in [4]), but statistical correlation (or significance).

Other related work includes correlation rule mining in [5]. It uses chi-square test to measure the correlation. A

correlation rule is a set of correlated items. It does not perform pruning or summarization. [10] introduces a technique for clustering association rules. It mainly deals with generation of numeric associations, i.e., how to join the adjacent intervals to produce more general rules and fewer rules. Clearly, both works are different from ours.

3. Chi-Square Test for Independence and Correlation

Chi-square test statistics (χ^2) is a widely used method for testing independence and/or correlation [17]. In our proposed technique, it is used in pruning as well as in finding direction setting rules. Below, we give an introduction to chi-square test.

Essentially, χ^2 test is based on the comparison of observed frequencies with the corresponding expected frequencies. The closer the observed frequencies are to the expected frequencies, the greater is the weight of evidence in favor of independence.

Example 3: In a loan application domain, we have 500 people who applied for loan in a bank. Out of the 500 people, 300 had a job and 200 did not have a job. 280 people were granted loan and 220 were not. We also know that 200 people who had a job were granted loan, which can also be expressed as an association rule:

Job = yes \rightarrow Loan = approved

[sup = 200/500, conf = 200/300]

This information gives us a 2x2 contingency table containing four cells (Figure 2). Note that the table has only 1 degree of freedom, which is sufficient for our work. See [17] or any standard statistics text for more details about degrees of freedom.

	Loan = approved	Loan = not-approved	Row Total:
Job = yes	200	100	300
Job = no	80	120	200
Column Total:	280	220	500

Figure 2. Contingency table for *Job* and *Loan*

Our question is "Is loan approval dependent on whether one has a job or not?" In dealing with this problem we set up the hypothesis that the two attributes are independent. We then compute the expected frequency for each cell as follows: Of the 500 people included in the sample, 300 (60% of the total) had a job, while 200 (40% of the total) had no job. If the two attributes are truly independent, we would expect the 280 approved cases to be divided between job = yes and job = no in the same ratio (60% and 40%); similarly, we would expect the 220 not-approved cases to be divided in the same fashion.

χ^2 is used to test the significance of the deviation from the expected values. Let f_0 be an observed frequency, and f be an expected frequency. The χ^2 value is defined as:

$$\chi^2 = \sum \frac{(f_0 - f)^2}{f}$$

	y	$\neg y$	Row Total
Satisfies X	Support count of $X \rightarrow y$	Support count of $X \rightarrow \neg y$	Support count of X in data
Does not satisfy X	Support count of $\neg X \rightarrow y$	Support count of $\neg X \rightarrow \neg y$	Support count of $\neg X$ in data
Column Total:	Support count of y in data	Support count of $\neg y$ in data	Total no. of cases in data

Figure 3. A generic 2x2 contingency table

A χ^2 value of 0 implies the attributes are statistically independent. If it is higher than a certain threshold value (e.g. 3.84 at the 95% significance level [17]), we reject the independence assumption. For our example, we obtain $\chi^2 = 34.63$. This value is much larger than 3.84. The independence assumption is rejected. Thus, we say that the loan approval is correlated to (or dependent on) whether one has a job.

However, we still do not know how they are correlated, e.g., whether having a job is *positively* or *negatively* correlated to the approval of a loan. Below, we give the definitions in the context of association rule mining.

Definition 1 (correlated): Let s be a minimum support and c be a significance level. X and y of a rule, $X \rightarrow y$, are said to be (s, c) correlated (*hereafter, merely correlated*) if the following two conditions are met:

1. The rule's support exceeds s .
2. The χ^2 value for the rule with respect to the whole data exceeds the χ^2 value at the significance level c .

Definition 2 (uncorrelated or independent): Let s be a minimum support and c be a significance level. X and y of a rule, $X \rightarrow y$, are said to be (s, c) uncorrelated (*hereafter, merely uncorrelated or independent*) if the following two conditions are met:

1. The rule's support exceeds s .
2. The χ^2 value for the rule with respect to the whole data does not exceed the χ^2 value at the significance level c .

Similar to those in [5], we define three types of correlation of a rule. For this work, we also call them the directions of a rule.

Definition 3 (type of correlation or direction):

Positive correlation: if X and y of a rule r , $X \rightarrow y$, are correlated and $f_o / f > 1$, we say that r is a *positive correlation* (or we say r is *significant*). We use 1 to denote a positive correlation. We also say that the *direction* of r is 1.

Negative correlation: if X and y of a rule r , $X \rightarrow y$, are correlated but $f_o / f < 1$, we say that r is a *negative correlation*. We use -1 to denote a negative correlation. We also say that the *direction* of r is -1.

Independence: if X and y of a rule r , $X \rightarrow y$, are independent, we say that r shows *independence*. We use 0 to denote independence. The *direction* of r is 0.

In this paper, we are only interested in the positively correlated rules (see “Additional features” in Section 6 if the other two types of rules are also of interest). The rule in Example 3 represents a positive correlation because its

observed frequency is 200 (i.e., the support count of the rule) and its expected frequency is only 168 ($= 300 * 280 / 500$).

In general, computing the type of correlation (or the direction) of an association rule r , $X \rightarrow y$, is to compare the rule with the whole population or the whole data set. Or more specifically, it is to compare with the rule that has the same conclusion as r but no condition, i.e., $\rightarrow y$. The generic contingency table used is shown in Figure 3.

Note that in the table, we use the term *support count* rather than frequency to conform to association rule mining. All support counts are available from the two rules (i.e., $X \rightarrow y$ and $\rightarrow y$), after running a mining algorithm. For example, to test the rule in Example 3 (R1), we compare it with R2:

R1: Job = yes \rightarrow Loan = approved
[sup = 200/500, conf = 200/300]

R2: \rightarrow Loan = approved
[sup = 280/500, conf = 280/500]

Clearly, these two rules specify completely the contingency table in Figure 2 (or Figure 3). R2, which has no condition, gives the *column total* for approved (280) and the total number of cases (500) in the data. R1, which represents the first cell in the table, also has the number 300 for the row total (which is simply the support count of Job = yes). With all this information, the rest of the cells can be computed.

4. Direction Setting Rules

Direction setting (DS) rules are the positively correlated association rules that set the directions for non-direction setting (non-DS) rules to follow. We give some definitions.

Definition 4 (direction setting rule): A rule r is a *direction setting* (DS) rule, if it satisfies the following conditions:

1. It has the direction of 1 (positive correlation).
2. Its direction is not an element of the set of *expected directions*.

Definition 5 (expected directions): The set of *expected directions* of a rule r is defined as follows:

1. If r is a 1-condition rule, the set of expected direction is $\{0\}$ (i.e. the condition and the consequent are expected to be independent in the absence of prior knowledge).
2. If r is a k -condition rule ($k > 1$) of the form:

$r: a_1, a_2, \dots, a_k \rightarrow y$

the set of expected directions is computed as follows: We view r as a combination of 2 rules, a 1-condition rule r_1 and a $(k-1)$ -condition rule r_{rest} ,

with the same consequent y .²

$r_1: a_i \rightarrow y \quad r_{rest}: a_1, a_2, \dots, a_j \rightarrow y$
 where $\{a_1, a_2, \dots, a_j\} = \{a_1, a_2, \dots, a_k\} - \{a_i\}$. The expected direction for this combination, denoted by E_i , is defined below ($x.dir$ denotes the direction of rule x):

- (a). if $r_1.dir = r_{rest}.dir = 1$ then $E_i = 1$.
- (b). if $(r_1.dir = 1 \text{ and } r_{rest}.dir = 0)$ or $(r_1.dir = 0 \text{ and } r_{rest}.dir = 1)$ then $E_i = 1$.
- (c). if $r_1.dir = r_{rest}.dir = 0$ then $E_i = 0$.
- (d). otherwise, $E_i = \text{unknown}$.

Since there are k such combinations, the set of expected directions of r is $\{E_i\}, i = 1, \dots, k$.

Lemma: All positively correlated 1-condition rules are direction setting rules.

Proof: It follows directly from Definition 4 and 5.

Notes about the above two definitions:

- In Definition 5, the assumptions (a) to (d) are reasonable because they follow our human intuitions:
 - (a) Two positive correlations generally lead to another positive correlation (Example 2).
 - (b) A positive correlation and an independence generally suggest a positive correlation.
 - (c) Two independence relations lead to another independence relation.
 - (d) For the rest of the situations, we either cannot decide what is the expected direction (e.g., 1 and -1 combination), or we simply do not care (e.g., -1 and -1 combination).
- The second condition (2) of Definition 4 basically says that if a rule's direction is expected with respect to any one combination, it is not a DS rule. This is reasonable because as long as there is one possible justification for r 's direction, we cannot say that r sets a new direction. Thus, r is not a DS rule.

To further explain the definitions, we list all possible direction combinations of r_1 , r_{rest} and r (Figure 4) using the notation:

$$r_1.dir, r_{rest}.dir := r.dir \text{ or } r_{rest}.dir, r_1.dir := r.dir$$

A (1) <u>1, 1 := 1</u>	B (1) <u>1, 0 := 1</u>	C (1) <u>0, 0 := 0</u>
(2) 1, 1 := -1	(2) 1, 0 := -1	(2) 0, 0 := -1
(3) 1, 1 := 0	(3) 1, 0 := 0	(3) 0, 0 := -1
D (1) -1, -1 := -1	E (1) -1, 0 := -1	F (1) -1, 1 := 0
(2) <u>-1, -1 := 1</u>	(2) <u>-1, 0 := 1</u>	(2) <u>-1, 1 := 1</u>
(3) -1, -1 := 0	(3) -1, 0 := 0	(3) -1, 1 := -1

Figure 4. All possible direction combinations of r_1 , r_{rest} and r

A(1), B(1) and C(1) conform to our expectations (a), (b) and (c) in Definition 5. That is, given the directions of r_1 and r_{rest} on the left-hand-side of “:=”, if r has the direction

on the right-hand-side of “:=”, we say r 's direction is expected. In fact, D(1), E(1) and F(1) can also be seen as expected. However, they are not interesting because we are only interested in the positively correlated rules. The interesting situations occur in C(2), D(2), E(2) and F(2) (shaded in Figure 4) because the direction of r is 1, but the expected direction of the r_1 and r_{rest} combination is 0 or unknown, which are our cases (c) and (d) in Definition 5. When such situations occur, we say that r sets a new direction. We call r a *potential* DS rule. r will be called a DS rule if for all possible combinations of r_1 and r_{rest} , the direction of r is different from the expected directions. Since we are only interested in positive correlations, the remaining situations in Figure 4 are not relevant (see also Section 6).

Definition 6 (non-direction setting rules): A *non-direction setting* (non-DS) rule is a positively correlated rule that is not a DS rule.

5. Pruning the Discovered Association Rules and Finding Direction Setting Rules

This section presents the basic ideas of pruning and finding direction setting rules. The next section gives the detailed algorithm, which performs both tasks.

5.1. Pruning of association rules

Section 3 shows that to test for correlation between the condition and the consequent of a rule $r, X \rightarrow y$, we compare it with the rule $R, \rightarrow y$, to see whether r is significant with respect to R . Those rules that are not positively correlated are removed. However, we can do much better than that. Consider again Example 1 (which is reproduced here as Example 4).

Example 4: We have the following two rules,

$R: \text{Job} = \text{yes} \rightarrow \text{Loan} = \text{approved}$
 [sup = 60%, conf = 90%]
 $r: \text{Job} = \text{yes}, \text{Credit_history} = \text{good} \rightarrow$
 $\text{Loan} = \text{approved}$ [sup = 40%, conf = 91%]

If we know R , then r is of limited use because it gives little extra information. Its slightly higher confidence is more likely due to chance than to true correlation.

We say r can be pruned with respect to R because within the subset of data cases covered by R , r is not significant. (A rule *covers* a set of data cases, if the data cases satisfy the conditions of the rule.) The pruning proceeds as that in Section 3. However, instead of using the whole data set, here we test the correlation of r with respect to R as r only covers a subset of the data cases that are covered by R . If r does not show a positive correlation with respect to R , it should be pruned. In general, pruning is done as follows:

- Given a rule r , we try to prune r using each *ancestor rule* R (which has the same consequent as r but fewer or 0 conditions) of r . That is, we perform a χ^2 test on r with respect to R . If the test shows a positive correlation, it is kept. Otherwise, r is pruned. The reason for the pruning is because within the data covered by R , r is not significant.

² We can also view the rule r as a combination of more than two shorter rules. The numbers of conditions in the shorter rules do not have to be 1 and $k-1$, but any possible partition. However, these alternatives require much more computation. Conceptually, they do not seem to have any advantage over our 1 and $k-1$ combination.

5.2 Finding DS rules

The process for finding DS rules is as follows: χ^2 test is first used to evaluate each 1-condition rule to determine its direction status, i.e., 1 (positive correlation), -1 (negative correlation), or 0 (independence). Then, it proceeds level-by-level to analyze each rule and decide whether it follows the direction that has already been set by rules at previous levels, or whether it sets a new direction. That is, at level 2 we analyze only 2-condition rules, at level 3, we only analyze 3-condition rules and so on. (For easy discussion, from now on we will use the level number and the number of conditions of a rule interchangeably). The analysis proceeds as follows: At level k ($k > 1$), for each k -condition rule r , we first determine its direction. We then examine each combination of 1-condition rule r_1 and $(k-1)$ -condition rule r_{rest} of r to determine whether r follows the expected direction set by r_1 and r_{rest} . If r follows the direction set by at least *one* such combination, we say r is not a DS rule. If r does not follow the direction set by any combination, we say r sets a new direction, and it is a DS rule.

The computation with each combination can be illustrated using the following if-statement (Figure 5) (the corresponding situations in Figure 4 are also indicated):

```

1  A      if  $r_1.dir = r_{rest}.dir = 1$  then
2    A(1)   if  $r.dir = 1$  then  $r$  is not a DS rule
3    A(2)(3) else  $r.dir = \text{undefined}$  3 /* see footnote */
4  B      elseif ( $r_1.dir = 1$  and  $r_{rest}.dir = 0$ ) or
           ( $r_1.dir = 0$  and  $r_{rest}.dir = 1$ ) then
5    B(1)   if  $r.dir = 1$  then  $r$  is not a DS rule
6    B(2)(3) else  $r.dir = \text{undefined}$  /* see footnote */
7  C      elseif  $r_1.dir = r_{rest}.dir = 0$  then
8    C(1)   if  $r.dir = 0$  then nothing
9    C(2)   elseif  $r.dir = 1$  then  $r$  is a potential DS rule
10   C(3)   else  $r.dir = \text{undefined}$  /* see footnote */
11 D, E, F else
12   D(2), E(2), F(2) if  $r.dir = 1$  then  $r$  is a potential DS rule
13   D(1)(3), E(1)(3), F(1)(3) else  $r.dir = \text{undefined}$  /* see footnote */
14 endif

```

Figure 5. Identifying a potential DS rule

We now state the key theorem of our technique.

Theorem: Using the above procedure to identify DS rules, every non-DS rule, $X \rightarrow y$, is a combination of one or more DS rules and zero or more 0 direction 1-condition rules with the same consequent y . That is, with a suitable re-arrangement of its conditions (as the ordering of the conditions are not important) any non-DS rule can be expressed as:

$$dsr + zr^*$$

where dsr is any DS rule, and zr is any 0 direction 1-condition rule.

Proof: See [14] for the full proof.

³ We need to assign *undefined* here to prevent some undesirable situations, see [14] for details.

Note that the above theorem does not say that any possible combination of some DS rules and/or independence rules (0 direction) is a non-DS rule. The reason is that such a combination may not meet the minimum support, or may not show positive correlation.

From the definition of DS rules and the theorem, we can derive three important points:

- Every DS rule r is *unexpected* with respect to all r_1 and r_{rest} combinations because r does not follow their directions.
- After seeing the DS rules, the directions of non-DS rules are no longer surprising as they are just some combinations of DS rules and independence rules (the direction of 0). That is, the non-DS rules are somewhat expected if the DS rules are known.
- DS rules can guide the user to see the related non-DS rules, if he/she is interested. The non-DS rules can provide further details with regard to the DS rules.

These points enable us to build a simple user interface that allows the user to focus on the essential aspects (DS rules) of the domain and selectively view the relevant details (non-DS rules). See Section 7.

6. The Algorithm

Figure 6 gives the algorithm (called P-DS) for both pruning and finding DS rules. The input parameters are F and T . F is the set of discovered large rules. T is the χ^2 value at a particular significance level. Two points to be noted:

- In the definitions of DS and non-DS rules, we did not mention how pruning is related. Clearly, those pruned rules will not be included in the set of DS rules or the set of non-DS rules.
- For easy understanding, the algorithm is presented as a post-processing method of the discovered rules. However, it can be easily incorporated into a rule miner itself. In fact, we implemented our system that way.

The algorithm processes the discovered rules level-by-level (line 1) from level-1 to level- n (where n is the highest level). For each rule r , $X \rightarrow y$, at a particular level the algorithm works as follows: The procedure *compDir* computes the type of correlation (or direction) of r (in line 2), which is given to $r.dir$. “ $\rightarrow y$ ” is a rule without any condition. In line 3 and 4, if r is a level-1 rule and its direction is 1, then r is a DS rule (DSR contains the set of all DS rules). If r ’s direction is not 1, we record that r is pruned by “ $\rightarrow y$ ” by assigning “ $\rightarrow y$ ” to $r.prune$ (line 5). This saved information is important for subsequent pruning (see procedure *evalPrune* in Figure 8). ($r.prune$ is initialized to 0, indicating that r is not pruned.) Line 6 processes r using all pairs of its ancestors r_1 and r_{rest} . In line 7, if r is pruned and r cannot be a DS rule, we can exit the for-loop. Anyone of the two conditions would not be sufficient for the exit. Evaluation of pruning is done in line 8 using the procedure *evalPrune*. Line 9 checks to see whether it has been determined that r cannot be a DS rule. If so, there is no need to proceed. From line 10-21, the algorithm analyzes r by considering the four cases. This part has been discussed in Section 5.2. $r.justify$ is used to

```

Procedure P-DS( $F, T$ )
1  for each  $r (X \rightarrow y) \in F$  from level-1 to level- $n$  do
2     $r.dir = compDir(r, "\rightarrow y", T)$ ;
3    if  $r$  is a level-1 rule then
4      if  $r.dir = 1$  then  $DSR = DSR \cup \{r\}$ ; /*  $r$  is a DS rule */
5      else  $r.prune = "\rightarrow y"$ ; /* record that  $r$  is pruned by " $\rightarrow y$ " */
6    else for each pair,  $r_1 (X_1 \rightarrow y)$  and  $r_{rest} (X_{rest} \rightarrow y)$  of  $r$ , where  $X_1 \in X$ , and  $X_{rest} = X - \{X_1\}$  do
7      if  $r$  is pruned and  $r$  cannot be a DS rule then exit-for;
8      if  $r$  is not pruned then  $evalPrune(r, r_{rest})$ ; /* no need to prune against  $r_1$  */
9      if  $r$  is still a potential DS rule then
10       if  $r_1.dir = r_{rest}.dir = 1$  then
11         if  $r.dir = 1$  then  $r.justify = \emptyset$  /*  $r$  is not a DS rule */
12         else  $r.dir = undefined$ 
13       elseif ( $r_1.dir = 1$  and  $r_{rest}.dir = 0$ ) or ( $r_1.dir = 0$  and  $r_{rest}.dir = 1$ ) then
14         if  $r.dir = 1$  then  $r.justify = \emptyset$  /*  $r$  is not a DS rule */
15         else  $r.dir = undefined$ 
16       elseif  $r_1.dir = r_{rest}.dir = 0$  then
17         if  $r.dir = 0$  then  $r.justify = \emptyset$  /*  $r$  is not a DS rule */
18         elseif  $r.dir = 1$  then  $r.justify = r.justify \cup \{(r_1, r_{rest})\}$  /*  $r$  is a potential DS rule */
19         else  $r.dir = undefined$ 
20       else if  $r.dir = 1$  then  $r.justify = r.justify \cup \{(r_1, r_{rest})\}$  /*  $r$  is a potential DS rule */
21       else  $r.dir = undefined$ 
22     endfor
23     if  $r.justify \neq \emptyset$  then /*  $r$  sets a new direction */
24       if  $r$  is pruned then  $r.dir = undefined$  /* Although  $r$  sets a new direction, it is pruned */
25       else  $DSR = DSR \cup \{r\}$  /*  $r$  is a DS rule */
26   endif
27 endfor;
28  $unprnRules = \{r \in F \mid r.dir = 1, r \text{ is not pruned}\}$ ; /* all the unpruned rules */
29  $non-DSR = unprnRules - DSR$  /* all the non-DS rules */

```

Figure 6. Algorithm P-DS

record all r_1 and r_{rest} combinations that justify r to be a potential DS rule (line 18 and 20). This information is helpful to the user in understanding the DS rules (see Section 7). When we know that r is not a DS rule (line 11, 14 and 17), no recording is needed (we set $r.justify = \emptyset$).

After completing the inner for-loop, if $r.justify \neq \emptyset$, r is justified to be a DS rule. However, if r can be pruned ($r.prune \neq 0$, see also Figure 8), then it will not be a DS rule (line 23). $r.dir$ is set to *undefined* (line 24) so that r can be used in the future to justify a rule for being a DS rule. If r cannot be pruned ($r.prune = 0$), then it is a true DS rule (line 25), and it is included in DSR . All unpruned rules are in $unprnRules$ (line 28), and all non-DS rules are in $non-DSR$ (line 29).

Below, we describe the two procedures *compDir* and *evalPrune*. Procedure *compDir* (Figure 7) uses the χ^2 test to compute the correlation or direction of r . R is an ancestor rule of r and both have the same consequent.

```

Procedure compDir( $r, R, T$ )
1  if  $\chi^2(r, R) > T$  then
2    if  $r.sup > r.cover * (R.sup / R.cover)$  then
3      return(1)
4    else return(-1)
5  else return(0)

```

Figure 7. Computing the direction or correlation of a rule r against a more general rule R

In line 1, if $\chi^2(r, R) > T$, we reject the independence assumption, i.e., the conditions and the consequent are correlated. Line 2-5 determine the type of correlation or

direction (Definition 3). " $r.cover * (R.sup / R.cover)$ " is the expected frequency. $x.cover$ is the number of data cases that satisfy the conditions of rule x .

The *evalPrune* procedure is shown in Figure 8. It tries to prune r using r_{rest} . In line 1, if r_{rest} itself has been pruned previously (i.e., $r_{rest}.prune \neq 0$). Then, the algorithm needs to find the rule that prunes r_{rest} . This is shown in line 2. In line 3 and 4, if r does not represent a positive correlation, it is pruned. We set $r.prune = r_{rest}$ to provide a link for pruning of higher level rules than r . In line 5, if r is a positive correlation, then we try to compare it with r_{rest} using chi-square test. If r does not show a positive correlation within the subset of the data covered by r_{rest} , it is pruned by r_{rest} (line 6).

```

Procedure evalPrune( $r, r_{rest}$ )
1  if  $r_{rest}.prune \neq 0$  then /* If  $r_{rest}$  has been pruned */
2     $r_{rest} = r_{rest}.prune$ ;
3  if  $r.dir \neq 1$  then
4     $r.prune = r_{rest}$ 
5  elseif  $compDir(r, r_{rest}) \neq 1$  then /*  $r$  is not significant compared to  $r_{rest}$  */
6     $r.prune = r_{rest}$  /*  $r$  can be pruned */

```

Figure 8. Evaluating pruning of a rule r with respect to its ancestor rule r_{rest}

Complexity of the algorithm: Let M be the number of discovered large rules, and N be the maximum number of conditions of a rule. The time complexity of the algorithm is thus $O(MN)$ (χ^2 test is very efficient and can be seen as constant). Since N is normally small, less

Table 1. Experiment results (minimum support = 1%; significance level = 95%)

1	2	3	4	5	6	7	8	9	10	
No.	Data sets	no. of Items	Target items	Large rules	PC. rules	PC rules Aft. pru	1-cond DS rules	Total no. DS rules	No. of conds in DS rules	Exe. time
1	Anneal	121	6	42893	18843	347	66	99	1.53	6.17
2	Austra	51	2	80000	40668	379	29	94	2.05	3.98
3.	Auto	98	7	80000	43781	2175	110	559	2.45	2.91
4	Breast-w	31	2	2757	2362	97	26	27	1.07	0.14
5	Chess	74	2	80000	31939	2314	50	253	2.60	2.40
6	Cleve	30	2	30435	11707	181	21	22	1.22	1.75
7	Crx	55	2	80000	40269	465	31	97	2.92	3.96
8	Diabetes	18	2	1196	565	44	14	14	1.00	0.05
9	German	66	2	80000	21386	496	32	138	2.44	3.44
10	Glass	21	7	1820	1036	153	35	42	1.18	0.07
11	Heart	23	2	10044	3055	131	18	18	1.00	0.50
12	Hepati	36	2	80000	23097	201	22	41	2.44	3.81
13	Horse	66	2	79965	24368	314	41	81	2.26	2.23
14	Hypo	50	2	80000	32543	299	23	49	1.96	3.92
15	Iono	314	2	34463	13710	682	238	405	1.49	1.04
16	Led7all	15	10	1692	1389	691	70	72	1.08	0.06
17	Lymph	51	4	80000	18324	511	40	89	1.89	3.26
18	Pima	18	2	1196	565	44	14	14	1.00	0.04
19	Sick	109	2	80000	23538	456	35	73	1.65	2.69
20	Sonar	82	2	80000	24005	1299	42	54	1.65	3.09
21	Tic-tac	28	2	8315	2712	290	14	90	2.71	0.31
22	Vehicle	68	4	80000	49538	3259	90	474	2.55	3.03
23	Waveform	211	3	35888	30971	1835	174	293	1.67	0.58
24	Wine	37	3	27143	15013	551	40	68	1.54	1.53
25	Zoo	37	7	80000	59270	1159	59	204	2.05	4.56
26	Disease	19	2	1622	462	35	13	13	1.00	0.07
27	Profile	373	5	80000	29362	1229	77	217	2.28	4.97
28	Quality	45	2	67841	29776	295	40	89	2.88	2.34
29	Results	278	2	54481	16426	995	116	151	1.40	1.43
30	Traffic	79	3	74610	5916	231	17	94	2.70	4.10
Average		84	3	50545	20553	705	53	131	1.91	2.14
minimum support = 2%; significance level = 95%:										
Average		84	3	44280	23885	612	50	93	1.75	2.12
minimum support = 1%; significance level = 90%:										
Average		84	3	50545	23297	957	57	162	2.08	2.18
minimum support = 2%; significance level = 90%:										
Average		84	3	44280	25978	810	53	109	1.86	2.08

Table 1 gave the results obtained using the significance level of 95% for the χ^2 test and minsup = 1%. The other results can be found in [14], but they are summarized here. Below, we explain each column in Table 1. The final row gives the average value for each column.

Column 1: It gives the name of the data set (the last 5 are our real-life data sets). The number of records (or cases) in these data sets range from a few hundreds to tens of thousands.

Column 2: It shows the number of items in each data set. Recall, an item is an attribute and a value (or an interval) pair, i.e., (*<attribute>*, *<value>*).

Column 3: It shows the number of items (or values) in the target attribute of each data set.

Column 4: It gives the number of large rules generated from each data set. We can see that the number of large rules generated from an association rule miner is huge for each data set. Almost half of the data sets cannot be completed even under the hard limit of 80,000.

Column 5: It gives the number of positively correlated (PC) rules found in each data set. The number is much smaller. However, on average, there are still more than 20,000 of them.

Column 6: It gives the number of positively correlated (PC) rules after pruning. The number is reduced drastically. More than 96% of the PC rules are pruned. However, on average, the number of unpruned rules is still large, e.g., 705. Such a large number of rules are still

hard to be analyzed by a human user.

Column 7: It gives the number of 1-condition DS rules.

Column 8: It gives the total number of DS rules found in each data set. The number of DS rules is manageable and can be analyzed manually. On average, only 131 rules are DS rules. Many of them are 1-condition DS rules.

Column 9: It gives the average number of conditions in the DS rules. We can see that the DS rules are mostly short rules.

Column 10: It gives the running time in second (running on 512MB Sparc 2) for each data set. This time includes rule generation (data reside on disk), pruning and finding direction setting rules. We can see that the proposed technique is very efficient. We could not log any time for pruning and finding DS rules alone due to the efficiency.

From the summary of the other results (below Table 1), we make the following observation:

- When the minimum support increases, fewer DS rules are produced. When the significance level is lowered to 90% for the χ^2 test, slightly more DS rules are produced. This is because pruning is less stringent. However, the numbers of DS rules are still manageable. The running times are similar because we have the hard limit.

9. Conclusion

Data mining is to find patterns or regularities to summarize the data. If it also produces a huge number of patterns, it will be of limited use because a human user does not have the ability to analyze these patterns. However, if such a huge number of patterns do exist in the data, it will not be appropriate to arbitrarily discard any of them or to generate only a small subset of them. It is much more desirable if we can summarize them. This paper proposes such a technique. This technique first prunes off those rules that contain little extra information as compared to their ancestors, and then identifies the direction setting rules to give a global picture of the underlying relationships in the domain. Although the number of discovered associations can be very large, experimental results and real-life applications have shown that the number of direction setting rules is typically very small and with very few conditions. They can be manually inspected by a human user without too much effort.

Finally, we believe that in general any data mining technique that may potentially produce a large number of patterns should provide a technique to summarize the findings or the generated patterns. In this way, the user will be able to obtain an overall picture of the domain without being overwhelmed by a large number of detailed patterns. From this summarized information, he/she can then find some interesting aspects to focus on. The proposed technique represents a major step towards this direction.

References

- [1] Aggarwal, C., and Yu, P. "Online generation of association rules." *ICDE-98*, 402-411.
- [2] Agrawal, R., Imielinski, T., Swami, A. "Mining association rules between sets of items in large databases." *SIGMOD-1993*, 1993, 207-216.
- [3] Agrawal, R. and Srikant, R. "Fast algorithms for mining association rules." *VLDB-94*.
- [4] Bayardo, R., Agrawal, R., and Gunopulos, D. "Constraint-based rule mining in large, dense databases." To appear in *ICDE-99*, 1999.
- [5] Brin, S. Motwani, R. and Silverstein, R. "Beyond market basket: generalizing association rules to correlations." *SIGMOD-97*, 1997, 265-276.
- [6] Fayyad, U. M. and Irani, K. B. "Multi-interval discretization of continuous-valued attributes for classification learning." *IJCAI-93*, 1022-1027.
- [7] Han, J. and Fu, Y. "Discovery of multiple-level association rules from large databases." *VLDB-95*.
- [8] Klemetinen, M., Mannila, H., Ronkainen, P., Toivonen, H., and Verkamo, A.I. "Finding interesting rules from large sets of discovered association rules." *CIKM-1994*.
- [9] Kohavi, R., John, G., Long, R., Manley, D., and Pfleger, K. "MLC++: a machine learning library in C++." *Tools with artificial intelligence*, 740-743, 1994.
- [10] Lent, B, Swami, A., and Widom, J. "Clustering association rules." *ICDE-97*, 1997.
- [11] Liu, B., and Hsu, W. "Post-analysis of learned rules." *AAAI-96*, 1996, pp. 828-834.
- [12] Liu, B., Hsu, W., and Chen, S. "Using general impressions to analyze discovered classification rules." *KDD-97*, 1997.
- [13] Liu, B., Hsu, W. and Ma, Y. "Integrating classification and association rule mining." *KDD-98*.
- [14] Liu, B., Hsu, W. and Ma, Y. "Summarizing the discovered associations using direction setting rules." Technical Report, SoC, 1999.
- [15] Mahta, M., Agrawal, R. and Rissanen, J. "SLIQ: A fast scalable classifier for data mining." *EDBT-96*.
- [16] Merz, C. J. and Murphy, P. UCI repository of machine learning databases, 1996.
[<http://www.cs.uci.edu/~mllearn/MLRepository.html>]
- [17] Mills, F. *Statistical Methods*, Pitman, 1955.
- [18] Ng, R. T. Lakshmanan, L. Han, J. "Exploratory mining and pruning optimizations of constrained association rules." *SIGMOD-98*, 1998.
- [19] Padmanabhan, B., and Tuzhilin, A. "A belief-driven method for discovering unexpected patterns." *KDD-98*.
- [20] Pazzani, M., Mani, S. and Shankle, W. R. "Beyond concise and colorful: learning intelligible rules." *KDD-97*, 1997.
- [21] Quinlan, R. *C4.5: program for machine learning*. Morgan Kaufmann, 1992.
- [22] Silberschatz, A., and Tuzhilin, A. "What makes patterns interesting in knowledge discovery systems." *IEEE Trans. on Know. and Data Eng.* 8(6), 1996, 970-974.
- [23] Silverstein, C., Brin, S., Motiwani, R., and Ullman, J. "Scalable techniques for mining causal structures." *VLDB-98*, 1998.
- [24] Srikant, R. and Agrawal, R. "Mining quantitative association rules in large relational tables." *SIGMOD-96*.
- [25] Srikant, R., Vu, Q. and Agrawal, R. "Mining association rules with item constraints." *KDD-97*, 1997, 67-73.
- [26] Toivonen, H. Klemetinen, M., Ronkainen, P., Hatonen, K., and Mannila, H. "Pruning and grouping discovered association rules." *Mlnet Workshop on Statistics, Machine Learning, and Discovery in Databases*, 1995, 47-52.