

# Mining Compressed Frequent-Pattern Sets\*

Dong Xin

Jiawei Han

Xifeng Yan

Hong Cheng

University of Illinois at Urbana-Champaign  
Urbana, IL 61801, USA

## Abstract

A major challenge in frequent-pattern mining is the **sheer** size of its mining results. In many cases, a high *min\_sup* **threshold** may discover only commonsense patterns but a low one may generate an explosive number of output patterns, which severely **restricts** its usage.

In this paper, we study the problem of compressing frequent-pattern sets. Typically, frequent patterns can be clustered with a **tightness** measure  $\delta$  (called  $\delta$ -cluster), and a *representative pattern* can be selected for each cluster. Unfortunately, finding a minimum set of representative patterns is NP-Hard. We develop two greedy methods, **RPglobal** and **RPlocal**. The former has the guaranteed compression bound but higher computational complexity. The latter sacrifices the theoretical bounds but is far more efficient. Our performance study shows that the compression quality using **RPlocal** is very close to **RPglobal**, and both can reduce the number of closed frequent patterns by almost two orders of **magnitude**. Furthermore, **RPlocal** mines even faster than **FPClose**[11], a very fast closed frequent-pattern mining method. We also show that **RPglobal** and **RPlocal** can be combined together to balance the quality and efficiency.

## 1 Introduction

Frequent-pattern (or itemsets) mining has been a focused research theme in data mining due to its broad applications at mining association [2, 3], **correlation** [6], **causality** [17], sequential patterns [4], episodes [15],

multi-dimensional patterns [14], max-patterns [5], partial periodicity [12], emerging patterns [8], and many other important data mining tasks.

The problem of frequent-itemsets mining can be defined as follows. Given a **transaction database**, let  $\mathcal{O} = \{o_1, o_2, \dots, o_d\}$  be the set of items that appear in the database,  $\mathcal{T} = \{t_1, t_2, \dots, t_k\}$  be the transaction set, and  $I(t_i) \in \mathcal{O}$  be the set of items in transaction  $t_i$ . For any itemset  $P$ , let  $O(P)$  be the corresponding set of items, and  $T(P) = \{t \in \mathcal{T} | O(P) \subseteq I(t)\}$  be the corresponding set of transactions. We say  $O(P)$  is the *expression* of  $P$ , and  $|T(P)|$  is the *support* of  $P$ . An itemset  $P$  is *frequent* if  $|T(P)| \geq \text{min\_sup}$ , where *min\_sup* is a user-specified threshold. The task of frequent-itemsets mining is to **find all the frequent itemsets**.

There have been many **scalable** methods developed for frequent-pattern mining [3, 13, 19]. However, the real bottleneck of the problem is not at the efficiency but at the **usability**. Typically, if *min\_sup* is high, mining may generate only commonsense patterns, however, with a low *min\_sup*, it may generate an explosive number of results. This has severely restricted the usage of frequent-pattern mining.

To solve this problem, it is natural to explore how to “compress” the patterns, i.e., find a **concise** and **succinct** representation that describes the whole collection of patterns. Two major approaches have been developed in this direction: lossless compression and lossy **approximation**. The former, represented by the *closed frequent itemsets* [16, 18, 19] (a frequent itemset  $P$  is *closed* if there is no itemset  $P'$  such that  $O(P) \subset O(P')$  and  $T(P) = T(P')$ ), emphasizes too much on the *supports* of patterns so that its compression power is quite limited. The latter, represented by the **maximal frequent itemsets** [5, 10, 11] (a frequent itemset  $P$  is *maximal* if there is no frequent itemset  $P'$  such that  $O(P) \subset O(P')$ ), as well as the **boundary cover sets** proposed recently [1], only consider the *expressions* of patterns, while the *support* information in most of the itemsets is lost. To achieve high-quality pattern compression, it is desirable to build up a pattern compression framework that concerns both the *expressions* and *supports* of the patterns. A motivation example is shown as follows.

**Example 1** Table 1 shows a subset of frequent item-

\* The work was supported in part by the U.S. National Science Foundation NSF IIS-02-09199/IIS-03-08215.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

sets on accidents data set [9], where 39, 38, 16, 18, 12, 17 are the name of individual items. The closed itemsets cannot get any compression on this subset. The maximal itemsets will only report the itemset  $P_3$ . However, we observe that itemsets  $P_2, P_3$  and  $P_4$  are significantly different w.r.t. their supports, and the maximal itemset totally loses this information. On the other hand, the two pairs  $(P_1, P_2)$  and  $(P_4, P_5)$  are very similar w.r.t. both expressions and supports. We suggest a high-quality compression as  $P_2, P_3$  and  $P_4$ .

Table 1: A Subset of Frequent Itemsets

ID	Itemsets	Support
$P_1$	{38, 16, 18, 12}	205227
$P_2$	{38, 16, 18, 12, 17}	205211
$P_3$	{39, 38, 16, 18, 12, 17}	101758
$P_4$	{39, 16, 18, 12, 17}	161563
$P_5$	{39, 16, 18, 12}	161576

A general proposal for high-quality compression is to cluster frequent patterns according to certain similarity measure, and then select and output only a *representative pattern* for each cluster. However, there are three **crucial** problems that need to be addressed: (1) how to measure the similarity of the patterns, (2) how to define quality guaranteed clusters where there is a representative pattern best describing the whole cluster, and (3) how to efficiently discover these clusters (and hence the representative patterns)?

This paper addresses these problems. First, we propose a distance measure between two frequent patterns, and show it is a valid distance **metric**. Second, we define a clustering **criterion**, with which, the distance between the representative pattern and every other pattern in the cluster is bounded by a threshold  $\delta$ . The objective of the clustering is to minimize the number of clusters (hence the number of representative patterns). Finally, we show the problem is **equivalent** to set-covering problem, and it is NP-hard w.r.t. the number of the frequent patterns to be compressed. We propose two greedy algorithms: the first one, **RPglobal**, has bounded compression quality but higher computational complexity; whereas the second one, **RPlocal**, sacrifices the theoretical bound but is far more efficient. Our performance study shows that the quality of the compression using **RPlocal** is very close to **RPglobal**, and both can reduce the number of patterns generated by about two orders of magnitude w.r.t. the original collection of closed patterns. Moreover, **RPlocal** directly mines representative patterns from database and runs even faster than **FPClose**[11], a fast closed frequent-itemset mining algorithm. We also show that **RPglobal** and **RPlocal** can be **integrated** together to balance the quality and efficiency.

The remaining of the paper is organized as follows. In Section 2, we formally introduce the problem. The NP-hardness is proved in Section 3. Section 4 proposes the **RPglobal** and **RPlocal** methods. Our performance study is presented in Section 5. A discussion on po-

tential extensions is in Section 6, and we conclude the study in Section 7.

## 2 Problem Statement

In this section, we first introduce a new distance measure on closed frequent patterns, and then discuss the clustering criterion.

### 2.1 Distance Measure

**Definition 1 (Distance measure)** Let  $P_1$  and  $P_2$  be two closed patterns. The distance of  $P_1$  and  $P_2$  is defined as:

$$D(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}$$

**Example 2** Let  $P_1$  and  $P_2$  be two patterns:  $T(P_1) = \{t_1, t_2, t_3, t_4, t_5\}$  and  $T(P_2) = \{t_1, t_2, t_3, t_4, t_6\}$ , where  $t_i$  is a transaction in the database. The distance between  $P_1$  and  $P_2$  is  $D(P_1, P_2) = 1 - \frac{4}{6} = \frac{1}{3}$ .

**Theorem 1** The distance measure  $D$  is a valid distance metric, such that:

1.  $D(P_1, P_2) > 0, \forall P_1 \neq P_2$
2.  $D(P_1, P_2) = 0, \forall P_1 = P_2$
3.  $D(P_1, P_2) = D(P_2, P_1)$
4.  $D(P_1, P_2) + D(P_2, P_3) \geq D(P_1, P_3), \forall P_1, P_2, P_3$

**Proof.** It is easy to verify that the first three properties are true. We prove the fourth statement is true.

To simplify the presentation, we define the following variables:  $|T(P_1)| = a$ ,  $|T(P_2)| = b$ ,  $|T(P_3)| = c$ ,  $|T(P_1) \cap T(P_2)| = b_1$ ,  $|T(P_2) - T(P_1) \cap T(P_2)| = b_2$ ,  $|T(P_1) \cap T(P_3)| = c_1$ ,  $|T(P_3) - T(P_1) \cap T(P_3)| = c_2$ ,  $|T(P_1) \cap T(P_2) \cap T(P_3)| = d_1$ , and  $|T(P_2) \cap T(P_3) - T(P_1) \cap T(P_2) \cap T(P_3)| = d_2$ . The meanings of the variables are shown in Fig. 1.

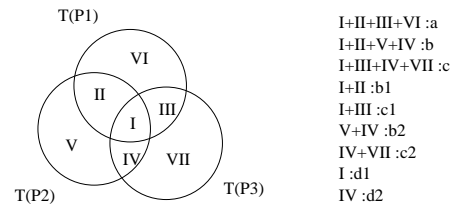


Figure 1: Meanings of Variables

Since  $(T(P_1) \cap T(P_2)) \cup (T(P_1) \cap T(P_3)) \subseteq T(P_1)$ , we have

$$\begin{aligned} & |T(P_1) \cap T(P_2)| + |T(P_1) \cap T(P_3)| \\ & - |T(P_1) \cap T(P_2) \cap T(P_3)| \leq |T(P_1)| \\ \implies & b_1 + c_1 - d_1 \leq a \end{aligned} \quad (1)$$

Plug in all the variables into the distance definition,

$$\begin{aligned} & D(P_1, P_2) + D(P_2, P_3) \geq D(P_1, P_3) \\ \iff & \frac{b_1}{a+b_2} + \frac{c_1}{a+c_2} \leq 1 + \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1-d_2} \end{aligned}$$

Using Eq. (1), we have:

$$\begin{aligned}
& 1 + \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1-d_2} \\
\geq & 1 + \frac{d_1}{b_1+b_2+c_1+c_2-d_1} \quad (d_2 \geq 0) \\
\geq & 1 + \frac{d_1}{a+b_2+c_2} \quad (Eq.1) \\
= & \frac{a+d_1+b_2+c_2}{a+b_2+c_2} \\
\geq & \frac{b_1+c_1+b_2+c_2}{a+b_2+c_2} \quad (Eq.1) \\
= & \frac{b_1+c_2}{a+b_2+c_2} + \frac{c_1+b_2}{a+b_2+c_2} \\
\geq & \frac{b_1}{a+b_2} + \frac{c_1}{a+c_2} \quad (a+b_2 \geq b_1, c_2 \geq 0) \\
& (a+c_2 \geq c_1, b_2 \geq 0)
\end{aligned}$$

Thus the fourth statement is true.  $\blacksquare$

*Remark.* The distance measure can be extended to the general frequent patterns **except that for non-closed patterns**, we may have  $D(P_1, P_2) = 0$  for some  $P_1 \neq P_2$ . This happens when  $P_1$  and  $P_2$  share the same support transactions set.

## 2.2 Clustering Criterion

By defining the distance on the set of transactions, the *support* information of patterns are well **incorporated**. We further consider the *expressions* of the patterns. Given two patterns  $A$  and  $B$ , we say  $B$  can be *expressed* by  $A$  if  $O(B) \subset O(A)$ . Following this definition, assume patterns  $P_1, P_2, \dots, P_k$  are in the same cluster. The representative pattern  $P_r$  of the cluster should be able to *express* all the other patterns. Clearly, we have  $\cup_{i=1}^k O(P_i) \subseteq O(P_r)$ .

Using the distance measure defined in Section 2.1, we can simply apply a clustering method, such as  $k$ -means, on the collection of frequent patterns. However, it introduces two problems. First, the quality of the clusters cannot be guaranteed; and second, it may not be able to find a representative pattern for each cluster (i.e., the pattern  $P_r$  may not belong to the same cluster). To overcome these problems, we introduce the concept of  $\delta$ -cluster, where  $\delta$  ( $0 \leq \delta \leq 1$ ) is the tightness measure of a cluster.

**Definition 2 ( $\delta$ -cluster)** A pattern  $P$  is  $\delta$ -covered by another pattern  $P'$  if  $O(P) \subseteq O(P')$  and  $D(P, P') \leq \delta$ . A set of patterns form a  $\delta$ -cluster if there exists a representative pattern  $P_r$  such that for each pattern  $P$  in the set,  $P$  is  $\delta$ -covered by  $P_r$ .

*Remark.* First, in  $\delta$ -cluster, one pattern can belong to multiple clusters. Second, using  $\delta$ -cluster, we only need to compute the distance between each pattern and the representative pattern in a cluster. Since a pattern  $P$  is  $\delta$ -covered by a representative pattern  $P_r$  only if  $O(P) \subseteq O(P_r)$ , we can simplify the distance calculation by only considering the supports of the patterns:  $D(P, P_r) = 1 - \frac{|T(P) \cap T(P_r)|}{|T(P) \cup T(P_r)|} = 1 - \frac{|T(P_r)|}{|T(P)|}$ . Finally, if we extend the distance definition to non-closed patterns, it is easy to verify that a non-closed

pattern must be  $\delta$ -covered by a representative pattern if its corresponding closed pattern is covered. We have the following **lemma** with the proof **omitted**.

**Lemma 1** Given a transaction database, a minimum support  $M$  and a cluster quality measure  $\delta$ , if a representative pattern set  $R$   $\delta$ -covers all the closed frequent patterns, then  $R$   $\delta$ -covers all the frequent patterns.

In the remaining of the paper, when we refer to frequent patterns, we mean *closed* frequent patterns. For simplicity, we use *cover* and *cluster* to represent  $\delta$ -cover and  $\delta$ -cluster, respectively.

If we restrict the representative pattern to be frequent, then the number of representative patterns (i.e., clusters) is no less than the number of maximal frequent patterns. This is because a maximal frequent pattern can only be covered by itself. In order to achieve more succinct compression, we relax the **constraints** on representative patterns, i.e., allow the supports of representative patterns to be *somewhat* less than  $\min\_sup$ ,  $M$ .

For any representative pattern  $P_r$ , assume its support is  $k$ . Since it has to *cover* at least one frequent pattern (i.e.,  $P$ ) whose support is at least  $M$ , we have

$$\delta \geq D(P, P_r) = 1 - \frac{|T(P_r)|}{|T(P)|} \geq 1 - \frac{k}{M}$$

That is,  $k \geq (1 - \delta)M$ . This is the  $\min\_sup$  for a representative pattern. To simplify the **notation**, we use  $\hat{M}$  to represent  $(1 - \delta)M$ .

The pattern compression problem is defined as follows.

**Definition 3 (Pattern Compression Problem)** Given a transaction database, a  $\min\_sup$   $M$  and the cluster quality measure  $\delta$ , the pattern compression problem is to **find a set of representative patterns  $R$ , such that for each frequent pattern  $P$  (w.r.t.  $M$ ), there is a representative pattern  $P_r \in R$  (w.r.t.  $\hat{M}$ ) which covers  $P$ , and the value of  $|R|$  is minimized.**

## 3 NP-Hardness

We show that the problem defined above is NP-Hard.

**Theorem 2** The problem of finding the minimum number of representative patterns is NP-hard.

**Proof.** We show that the pattern compression problem can be reduced from the set-covering problem.

First, for any pattern compression problem, we can construct a corresponding set-covering problem. There are two  $\min\_sup$ s in the pattern compression problem:  $\min\_sup$   $M$  and the representative pattern's  $\min\_sup$   $\hat{M}$ . We denote the set of frequent patterns (w.r.t.  $M$ ) as  $FP(M)$ , the set of frequent patterns (w.r.t.  $\hat{M}$ ) as  $FP(\hat{M})$ . For each pattern  $\hat{P} \in FP(\hat{M})$ , we **generate** a set whose elements are all the patterns  $P \in FP(M)$  which are *covered* by  $\hat{P}$ . The set-covering

problem is to find a minimum number of sets which cover all the elements, where each set corresponds to a representative pattern.

We then show that for any set-covering problem, we can construct a corresponding pattern compression problem. Let the given set-covering problem contain  $N$  elements  $(e_1, e_2, \dots, e_N)$  and  $K$  sets  $(S_1, S_2, \dots, S_K)$ . Each set  $S_i$  contains  $n_i$  elements  $(e_i^1, e_i^2, \dots, e_i^{n_i})$ . We assume that (1) there exist no two sets  $S_i$  and  $S_j$  such that  $S_i \subseteq S_j$ ; and (2) there exists no single set covering all the elements.

By considering these elements as individual items, we first construct  $\alpha (\geq 0)$  transactions, each of which contains  $N$  items  $(e_1, e_2, \dots, e_N)$ ; then construct  $\beta (\geq 1)$  transactions for each individual set (i.e., for each set  $S_i$ , there will be  $\beta$  transactions  $(e_i^1, e_i^2, \dots, e_i^{n_i})$ ). Now we have a database containing  $\alpha + \beta K$  transactions. If there is any individual item  $e_i$  whose support is less than  $\alpha + \beta K$ , we further insert transactions with only one item  $(e_i)$  until its support reaches  $\alpha + \beta K$ . Let  $M = \alpha + \beta K$  and  $\delta = \frac{\beta K - 1}{\alpha + \beta K}$ , then  $\hat{M} = \alpha + 1$ . Since  $\alpha$  and  $\beta$  can be chosen arbitrarily, the value of  $\delta$  can be selected as any value in  $(0, 1)$ .

Now we have a database where the longest pattern  $(e_1, e_2, \dots, e_N)$  is not frequent w.r.t.  $\hat{M}$ . It will not be considered as a representative pattern. Each set in the original set-covering problem corresponds to an itemset whose support is at least  $\alpha + \beta \geq \hat{M}$  (we denote the set of all of these itemsets as  $RP$ ) and can be considered as representative patterns. Each element in the set-covering problem corresponds to an item whose support is exactly  $M$  and has to be covered. We show that the solution of the pattern compression problem will only choose representative patterns from  $RP$ . This is because  $RP$  is the maximal pattern set w.r.t.  $\hat{M}$ . If there is a representative pattern  $P \notin RP$ , we can always find a  $P' \in RP$ , such that  $O(P) \subset O(P')$ . Since all the frequent patterns covered by  $P$  have supports at most  $M$ , they can also be covered by  $P'$ . We conclude that the optimal selection of the representative patterns corresponds to the solution of the original set-covering problem. ■

In the rest of the paper, we treat the following terms as equivalent: *element* vs. *frequent pattern* (w.r.t.  $M$ ); *set* vs. *frequent pattern* (w.r.t.  $\hat{M}$ ); and *set-cover* vs. *set of representative patterns*. For any frequent pattern  $P$  (w.r.t.  $\hat{M}$ ), we denote the set of patterns which can be covered by  $P$  as  $set(P)$ .

## 4 Discovering Representative Patterns

In this section, we describe algorithms for computing representative patterns.

### 4.1 The RPglobal Method

Generally, the size of the frequent patterns is quite large. It is undesirable to enumerate all the combinations to find the optimal selection. Since the problem

is equivalent to the set-covering problem, it is natural to consider some approximate algorithms available in the set-covering problem. The well-known one is the greedy algorithm [7] which iteratively finds the current largest set. The pseudo-code for the pattern compression problem is shown in Fig. 2. Since the precondition for this method is to collect the complete information over the elements and sets, we refer it as *global method* (in contrast to the *local method* to be discussed in the next section).

**Algorithm 1** (RPglobal) Compute Representative Patterns by Greedy Set-Covering.

**Input:** (1) A collection of frequent patterns  $FP$  w.r.t.  $\hat{M}$ , (2) a minimum support,  $M$ , and (3) a quality measure for clustering,  $\delta$ .

**Output:** The set of representative patterns.

**Method:** The algorithm is described in Figure 2.

```

BEGIN
  for each  $P \in FP$  s.t.  $support(P) \geq M$ 
    Insert  $P$  into the set  $E$ ;
  for each  $Q \in FP$ , s.t.  $Q$  covers  $P$ 
    Insert  $P$  into  $set(Q)$ ;
  while  $E \neq \emptyset$ 
    Find a  $RP$  that maximizes  $|set(RP)|$ ;
    for each  $Q \in set(RP)$ 
      Remove  $Q$  from  $E$  and the remaining sets;
  Output  $RP$ ;
END

```

Figure 2: The RPglobal algorithm

The code is self-explanatory. Following the result of greedy set-covering [7], the ratio between the number of the representative patterns selected by RPglobal and that of the optimal one is bounded.

**Theorem 3** Given a collection of frequent patterns  $\mathcal{F}$ , let the set of representative patterns selected by RPglobal be  $C_g$ , the set of optimal (i.e., minimal number) representative patterns be  $C^*$ , then  $|C_g| \leq |C^*| \times H(\max_{P \in \mathcal{F}} |set(P)|)$ , where  $H(n) = \sum_{k=1}^n \frac{1}{k}$ .

**Proof.** See [7].

The RPglobal method contains two steps. The first one is to collect the complete coverage information (i.e., find all the frequent patterns  $Q$  that can cover  $P$ ), and the second one is to find the set-covers (i.e., find the set of representative patterns). The greedy set-covering step can be implemented in time complexity of  $O(\sum_{P \in \mathcal{F}} |set(P)|)$  [7]. The computational challenge comes from finding the pattern coverage information. Note this coverage problem is different from closedness checking, which can be handled more efficiently because of the following reasons. First, closedness checking only needs to find one super-pattern



which **subsumes** the query pattern, **whereas** the coverage checking has to find *all* super patterns that can cover it. Second, the closedness checking can **utilize** transaction ID-based hash functions to do fast checking [19], while the coverage checking cannot benefit from it since there is a  $\delta$  tolerance between the support transaction sets. To **facilitate** the coverage search, we use an FP-tree-like structure [13] to index all the frequent patterns (w.r.t  $\hat{M}$ ). An example of FP-tree is shown in Fig. 3. The FP-tree has a head table associated with it. Single items are stored in the head table. The **entry** for an item also contains the head of a list that links all the nodes with the same name.

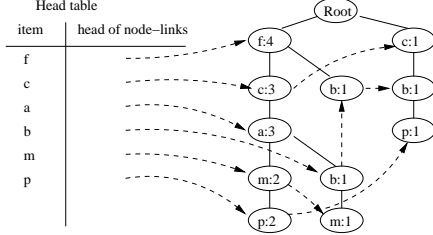


Figure 3: A sample FP-tree

The construction of the index tree is similar to FP-tree, except that in FP-tree, the counts of the nodes are updated by summing the counts of the inserted itemsets, while here, the counts of the nodes are updated by **choosing the maximum count over the inserted itemsets**. To differentiate from traditional FP-tree, we call our index tree as **RP-tree** (representative pattern tree). The coverage checking using RP-tree works as follows. Suppose the current pattern is  $Q$ ,  $O(Q) = \{o_1, o_2, \dots, o_k\}$  (items are ordered as in the RP-tree head table), and its support is  $C$ . The support region for a valid representative pattern is  $[C \times (1 - \delta), C]$ . Following the linked list of  $o_k$  in RP-tree, for each node  $n$  in the list, we test whether (1) the count is within the support region; and (2) the query itemset is a subset of the ancestors of  $n$ .

The worst case computation complexity for coverage checking could be  $O(|\mathcal{F}|^2)$ . The RPglobal method works well when  $|\mathcal{F}|$  is not large. However, when the number of frequent patterns to be compressed increases, the method does not scale well. It is necessary to develop an alternative method which discovers the set of representative patterns efficiently, while still preserves the high quality of the results.

## 4.2 The RPlocal Method

In this subsection, we introduce the idea of a *local method* and show how this method can be efficiently **incorporated** into the frequent-pattern mining process.

### 4.2.1 Local Greedy Method

Computing the complete coverage information is necessary for RPglobal, since the method needs to find a globally maximum set at each step. To develop

a **scalable** method, this expensive computational requirement has to be relaxed. Our objective is to report the representative patterns by an almost linear scan of the whole collection of patterns, without knowing the complete coverage information. The **intrinsic** relationship among the nearby patterns (according to the order generated by frequent pattern mining algorithms) can be utilized for this purpose.

Most frequent pattern mining algorithms conduct depth-first **enumerations** in the pattern space. It starts from an *empty* pattern set, recursively calls the pattern-growth routine to expand the pattern set. Since the individual items are sorted, at any stage of the algorithm, all the single items can be **partitioned** into three disjoint sets: the *conditional set* (the items appearing in the current pattern), the *todo-set* (the items to be expanded based on the current pattern) and the *done-set* (all the other items).

**Example 3** Fig. 4 shows a search space with five single items  $a, b, c, d, e$ . At the time when the depth-first search reaches pattern  $\{a, c\}$ , the *conditional set* is  $\{a, c\}$ , the *todo-set* is  $\{d, e\}$  and the *done-set* is  $\{b\}$ .

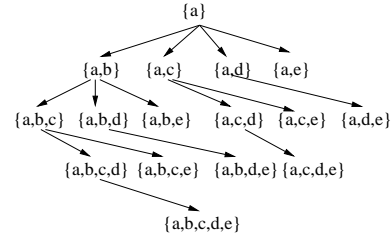


Figure 4: Depth-First Search in Pattern Space

The depth-first search scans each pattern twice: the first visit from its parent, and the second visit after finishing the calls to its children. One can **verify** that after a pattern is visited in its second time, all the patterns that can possibly cover it have been enumerated. The future patterns are not able to cover it.

We output a pattern in its second visit. The *local greedy method* sequentially scans the output patterns, at any time when an uncovered pattern (called **probe pattern**) is found, the algorithm finds the current largest set (i.e., a representative pattern) which covers it. Here the current largest set has the same meaning as it has in the global greedy method (i.e., the already covered pattern does not count for the set size). The following **theorem** shows a bound of the local method.

**Theorem 4** Given a collection of frequent patterns  $\mathcal{F}$ , let the set of representative patterns selected by the local method be  $C_l$ , the set of optimal representative patterns be  $C^*$ . Assume the minimum number of patterns that cover all the probe patterns be  $T$ . Then  $|C_l| < |C^*| \times (\sqrt{2T \times \max_{P \in \mathcal{F}} |\text{set}(P)|} + 1)$ .

**Proof.** See Appendix.

The difference between the local method and the global method is the selections of the **probe** patterns.

Clearly, if the probe patterns are selected as patterns in the current largest set, the local method is **identical** to the global method. Since the complete coverage information is not available, the bound on the local method is worse than the global method. However, in our experiments, we found that the performance of the local method is very close to that of the global method. This is because in the pattern compression problem, the layout of the patterns and their coverage is not **arbitrary**. Instead, most frequent patterns are strongly connected if they are in the same cluster, i.e., a pattern  $P$  is covered by the representative pattern  $P_r$  if and only if  $P_r$  subsumes  $P$  and the distance between  $P$  and  $P_r$  is within  $\delta$ . As a result, for each pattern  $P_r$ ,  $set(P_r)$  preserves local compactness, that makes the selections of probe patterns not a dominant factor for compression quality. Meanwhile, in most pattern compression problem, the sizes of sets are somewhat balanced. It is unlikely to have a very large set which is selected by the global method but missed by the local method, thus leads to a significant performance difference.

The local method relaxes the requirement of global comparison of sets sizes. However, it still needs to find the current largest set at each step, which involves the coverage checking for the future patterns. We further relax the method by finding a *reasonably large* set, instead of the *largest* set. The reasonably large set is expected to cover more future patterns. **Intuitively**, the **candidate** set should be as long as possible, since longer patterns generally have larger coverage. The candidate set should also contain more items within the *probe pattern's todo-set*. This is because items in todo-set are expected to appear more in the future.

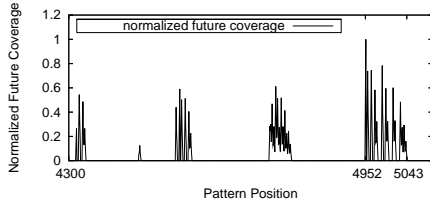


Figure 5: Patterns' Positions and Their Coverage

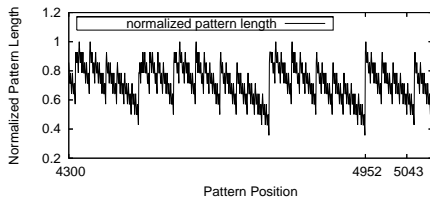


Figure 6: Patterns' Positions and Their Length

These **intuitions** are well justified by the real experimental statistics. Fig. 5 (*connect* data set [9] with  $minsup = 0.8 \times \#transactions$ ,  $\delta = 0.1$ ) shows the future coverage w.r.t. a probe pattern which is output at position 5043. The future coverage counts the number

of coverable patterns which are output after the probe pattern. We ignore the patterns which cannot cover the probe pattern (i.e., the future coverage is 0). The values on the  $y$ -axis are normalized w.r.t. the largest coverage. The  $x$ -axis is the order in which patterns are output. The probe pattern is first visited at position 4952. Fig. 6 shows the corresponding pattern length (normalized w.r.t. the longest length). We observe that the largest future coverage appears between the first visit and second visit of probe pattern, and it also has the longest length within the same region. Based on the above observations, we select the *reasonably large* set as the longest pattern, which can cover the probe pattern, among all the patterns between the first visit and second visit of the probe pattern (i.e., patterns expanded by the probe pattern).

Since the local method only requires the knowledge on already-discovered patterns, we further integrate it into frequent pattern mining process in order to improve the computational efficiency. The new method is called RPlocal.

#### 4.2.2 The Algorithm

We develop an FP-growth-like algorithm [13, 11] to discover representative patterns. Since FP-growth is a well-known method, we omit the detailed description here due to the limited space.

The RPlocal algorithm is described as below.

**Algorithm 2** (*RPlocal*) *Compute Representative Patterns by Local Search.*

**Input:** (1) a transaction database  $D$ , (2) a minimum support,  $M$  and (3) a quality measure for clustering,  $\delta$ .

**Output:** The set of representative patterns.

**Method:** The algorithm is described in Figure 7.

We explain the algorithm line by line. Line 1 picks an item to be expanded from the head table. Line 2 pushes the item onto a global stack  $IS$  which keeps track of the itemsets along the path from the root in the pattern space. Each entry of  $IS$  has the following fields: *item name*, *counts*, *covered*, and *cover pattern* (among its children). Line 3 checks whether closed **pruning** can be applied on  $IS(top)$ . We discuss the challenge and our solution to the closed pruning problem in the next subsection. Line 4 traces all the itemsets in the  $IS$  stack to check whether the current itemset can cover the itemsets in the stack. If yes, and the current pattern is longer than the one stored in the *cover pattern* field, then the *cover pattern* field is updated by the current pattern. Line 5 checks current itemset's support. If it is less than  $M$ , then it is not required to be covered. It also checks whether it is covered by a previous representative pattern (the previous representative patterns are indexed in RP-tree  $R$ , as shown in RPlglobal). Line 8-11 are the same as the FP-growth algorithm. It collects the todo-set based current itemset, constructs a conditional FP-tree and a

```

BEGIN
   $IS = \phi$ ; // global stack to keep itemsets
  scan  $D$ , create FP-tree  $F$ ;
  Initiate an empty RP-tree  $R$ ;
  call  $FPrepresentative(F, R)$ ;
END

procedure  $FPrepresentative(F, R)$  {
1. for each item in  $F.headtable$  {
2.   push item into  $IS$ ;
3.   if ( $closed\_pruning(IS(top)) = true$ ) continue;
4.    $set\_representative()$ ;
5.   if  $F.headtable[item].count < M$ 
       or  $coverage\_checking(IS(top), R) = false$ ;
6.      $IS[top].covered = true$ ;
   else
7.      $IS[top].covered = false$ ;
8.    $Todo - set = \{\text{frequent (w.r.t } \hat{M}) \text{ items based on}$ 
        $\text{the current conditional set}\}$ ;
9.   Build a new FP-tree  $F_{new}$  based on  $Todo - set$ ;
10.  Initiate  $F_{new}$ 's RP-tree,  $R_{new}$ ;
11.  call  $FPrepresentative(F_{new}, R_{new})$ ;
12.  if ( $(RP = get\_representative()) \neq NULL$ )
13.    Insert  $RP$  into  $R$  and its predecessor RP-trees;
14.    Output  $RP$ ;
15.  pop item from  $IS$ ;
  }
}

```

Figure 7: The RPlocal algorithm

conditional RP-tree for the coverage checking, then recursively calls the  $FPrepresentative$  routine with the new trees. The conditional RP-tree shares the same conditional set with the corresponding FP-tree. The approach is shown more efficient in [11] since the conditional RP-trees are more compact than the global RP-tree. Line 12 checks whether the current itemset can be a probe pattern. If it is, then the *cover pattern* stored in the  $IS$  stack is selected as a new *representative pattern*. Meanwhile, the *covered* fields of all the other itemsets in the  $IS$  stack are set as true if they can be covered by the new representative pattern. The new representative pattern is inserted into all RP-trees for future coverage checking.

#### 4.2.3 Prune Non-Closed Patterns

Here we discuss the implementation of *closed pruning* in the RPlocal algorithm.

Assume a pattern  $P$  is not closed, and the related closed pattern is  $P_c$ . There are two possibilities for  $P_c$ . One is  $(O(P_c) - O(P)) \cap done - set \neq \phi$ , then  $P_c$  is a pattern discovered before the first visit of  $P$ . The other is  $(O(P_c) - O(P)) \cap done - set = \phi$ , then  $P_c$  is a pattern expanded by  $P$  (i.e.,  $P_c$  is discovered between the first visit and the second visit of  $P$ ).

The second case is *intrinsically* handled by the RPlocal algorithm at line 8, where items with the same frequency as  $P$  are directly merged into the conditional

set, and the non-closed itemsets are skipped. The first case is more interesting w.r.t. the computation pruning. The following lemma has been widely used in most of the closed frequent pattern mining methods [18, 19], and we state it without proof.

**Lemma 2** *In the RPlocal method, for any pattern  $P$ , if there exists a pattern  $P_c$  which was discovered before the first visit of  $P$ , s.t.  $O(P) \subset O(P_c)$  and  $|T(P)| = |T(P_c)|$ , then all the patterns being expanded by  $P$  (i.e., patterns within the first and second visits of  $P$ ) are not closed.*

We call this pruning technique as *closed pruning*. The function *closed pruning* is to check whether the pattern is not closed w.r.t. a previously discovered pattern.

The challenge for closed pruning in the RPlocal algorithm is that only representative patterns are kept, and generally it is a small subset of closed patterns. It is not possible to check the closedness of a pattern using the previous outputs. Keeping all the closed patterns is one option. However, an interesting observation from our experiments shows that even without closed pruning, RPlocal runs faster than the closed frequent pattern mining algorithm. This is because the *coverage checking* in RPlocal is much more efficient than the *closedness checking* in closed frequent pattern mining since the number of representative patterns to be checked with is significantly less than the number of closed frequent patterns in closedness checking. Keeping all the closed patterns obviously will degrade the performance of RPlocal.

Instead of checking the closedness with the previous output, RPlocal uses a closedness checking method which tries to memorize the information of items in *done-set*. Since we only need to know whether an item is present or not in the closed pattern, we use 1 bit to represent an item's presence. If there are  $N$  single frequent items, we use a  $N$ -bit array (referred as *closed\_index*) for each pattern. The *closed\_index* of a pattern is computed by the bit-and operation between the *closed\_indices* of all the related transactions. An example on how to use *closed\_index* is shown as follows.

**Example 4** *Given a database having 5 transactions:  $\{f, c, a, m, p\}$ ,  $\{f, c, a, b, m\}$ ,  $\{f, b\}$ ,  $\{c, b, p\}$ ,  $\{f, c, a, m, p\}$ , we use  $N = 6$  bits for the *closed\_index*. Items  $f, c, a, b, m, p$  are assigned to the 1st to 6th bits, according to the computation order. The *closed\_indices* of transactions 1 and 5 are 111011, and the *closed\_index* of transaction 2 is 111110. The *closed\_index* of pattern  $\{c, a\}$  is 111010. Since item  $f$  is in the *done-set* of pattern  $\{c, a\}$  and  $f$ 's bit is 1. We conclude that the *closed pruning* can be applied on pattern  $\{c, a\}$ .*

To efficiently compute the *closed\_index* for each pattern, we attach *closed\_index* to each node in the FP-tree. The *closed\_index* can be aggregated along with the count measure, except that the count is updated by sum, while the *closed\_index* is updated by

bit-and. Since the *closed\_index* is attached to every node, the method is limited by the memory requirement. Fortunately, our task is not to identify all closed pruning. Instead, we aim to prune as much as possible, and the unpruned non-closed patterns will go to the coverage checking. The problem turns out: Given a fixed number of  $k$  for *closed\_index*, if the total number of single frequent items is larger than  $k$ , how to select  $k$  items from them, and how much pruning can be achieved?

It is natural to choose the first  $k$  items according to the computation order because the closed pruning checks patterns with items in *done-set*. The experimental statistics on *pumsb\_star* [9] data set is shown in Fig. 8, where we collect the percentage of closed pruning achieved, by setting  $k$  as 32 and 64. We observe this simple optimization works quite well. With only one integer ( $k = 32$ ) as *closed\_index*, the method misses less than 2% and 15% closed pruning when the number of frequent items are 2 and 6 times of  $k$ , respectively. Using two integers ( $k = 64$ ) as *closed\_index*, the method misses less than 1% of the closed pruning.

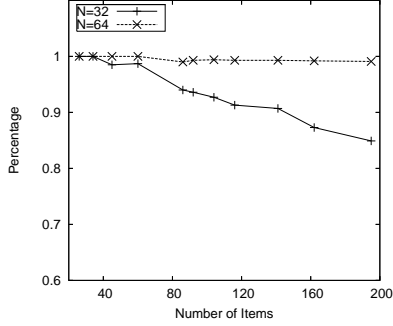


Figure 8: **Percentage of Closed Pruning achieved**

The similar phenomena are also observed on all the other data sets (e.g., *mushroom*, *connect*, *accidents*, *chess*) in our experimental evaluations. It is interesting to see that these limited  $k$  bits achieve good pruning percentages. We give a detailed explanation in the rest of this subsection.

Assume there are totally  $n$  independent frequent items, whose computation order is  $o_1, o_2, \dots, o_k, \dots, o_n$  (for simplicity, we assume the order of items keeps unchanged). We leave the first  $k$  for *closed\_index*, and  $r = n - k$  as left items. The percentage of the closed pruning by *closed\_index* is defined as function  $h(k, r)$ .

For any pattern, let the conditional set be  $\{o_{i_1}, o_{i_2}, \dots, o_{i_m}\}$ , where  $i_1 < i_2 < \dots < i_m$ , we say  $m$  is the *length* of the pattern and  $i_m$  is the *span* of the pattern. The position  $j$  is called a *hole* if  $j < i_m$  and  $j \notin \{i_1, i_2, \dots, i_m\}$ . If the set of holes is not empty (i.e., the *done-set* is not empty), then this pattern is possible to be subsumed by a previously output pattern (i.e., the closed pruning is possible to be applied on). A hole is *active* if the closed pruning takes effect.

The items in the conditional set are distributed into

two parts: the first  $k$  items set and the rest  $r$  items set. Let the number of items falling in the rest  $r$  items set be  $v$ , and the number of holes falling in the rest  $r$  items set be  $u$  (referred as  $(u, v)$ -configuration). To estimate the percentage of closed pruning for a  $(u, v)$ -configuration (defined as  $g(k, u, v)$ ), we need to further define two parameters: the expect number of active holes  $c$  and the maximal pattern length  $l$ .

Assume items are independent, every *hole* has an equal opportunity to be *active*. If there is *one* hole, which exists in the first  $k$  items, then the closed pruning is caught by the *closed\_index*, otherwise, it misses. Since there are at most  $l - v$  items falling into the first  $k$  items set, for each  $0 \leq i \leq m = \max(l - v, k)$ , there are  $\binom{k}{i}$  different patterns. For each pattern, the number of all different placements for  $c$  active holes is  $\binom{k-i+u}{c}$ , and the number of placements that all  $c$  active holes falling in the rest  $r$  items set is  $\binom{u}{c}$ . Thus the pruning percentage by *closed\_index* is  $1 - \frac{\binom{u}{c}}{\binom{k-i+u}{c}}$ . We have:

$$g(k, u, v) = \frac{\sum_{i=0}^m (1 - \frac{\binom{u}{c}}{\binom{k-i+u}{c}}) \times \binom{k}{i}}{\sum_{i=0}^m \binom{k}{i}}$$

Now we examine the value of  $h(k, r)$ . Among all patterns, there are two cases which are divided evenly. First, the last item is not in the conditional set. In this case, the pruning percentage is same as  $h(k, r - 1)$ . Second, the last item is in the conditional set. In this case, we enumerate all possible  $(u, v)$ -configurations. There are at most  $l - 1$  bits to be placed within the latter  $r$  items (since the last item is already in, there are  $r - 1$  selections). For each  $0 \leq i \leq m = \max(l - 1, r - 1)$ , there are  $\binom{r-1}{i}$  different  $(u, v)$  configurations ( $u = r - 1 - i, v = i + 1$ ), and for each configuration, the pruning percentage is  $g(k, r - 1 - i, i + 1)$ , we have:

$$h(k, r) = \frac{1}{2} h(k, r - 1) + \frac{\sum_{i=0}^m \binom{r-1}{i} \times g(k, r - 1 - i, i + 1)}{2 \sum_{i=0}^m \binom{r-1}{i}}$$

The base case is  $h(k, 0) = 1$ . We run simulations by setting  $k = 32$  and varying  $r$  from 0 to 64. We observe that, in most cases, the maximal pattern length is approximately proportional to the number of frequent items. Typically, we select the ratio as  $\frac{1}{3}$ , which is close to the experiments in *pumsb\_star* data set. The value of expected active closed bits  $c$  is varying from 1 to 4.

The simulation result in Fig. 9 shows that the percentage of pruning increases as  $c$  increases. This is because when  $c$  is large, the probability that at least one active holes are caught by *closed\_index* is high. Typically, when  $c = 4$ , the simulation curve is close to the real experimental results. Note  $c = 1$  is the base line where the corresponding curve represents the lower bound of closed pruning. We believe the



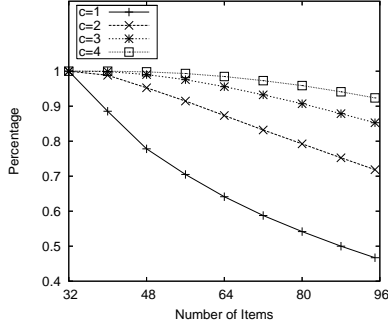


Figure 9: **Simulation Results: Percentage of Closed Pruning achieved**

*closed\_index* method has practical usage in real applications. The reason is as follows. The number of frequent patterns is exponentially explosive w.r.t. the number of items. Within the current computational limit, if the number of items is large, either the maximal length of the pattern is small or the pattern length is not small but the value of  $c$  is reasonably large (thus the output closed patterns can be kept in a reasonable size). In the former case, the effect of closed pruning (using whatever methods) to the whole computational efficiency is limited; while in the latter case, our study shows that *closed\_index* can achieve considerable improvement. Furthermore, the *closed\_index* approach only involves bit operations, which is very efficient.

### 4.3 Combining RPglobal and RPlocal

We have developed two algorithms for the pattern compression problem. The RPglobal method has guaranteed compression bound but is worse on scalability, whereas the RPlocal method is efficient but worse on the compression. In this subsection, we discuss a combined method: RPcombine.

The main idea of RPcombine is to first use RPlocal to get a small subset of candidate representative patterns, then use RPglobal to find the final results. To ensure that all the frequent patterns (w.r.t.  $M$ ) are  $\delta$ -covered, we need to choose the parameters for each step carefully.

Assume the quality measures for RPlocal and RPglobal are  $\delta_l$  and  $\delta_g$ , respectively. Any frequent pattern  $P$  must be  $\delta_l$ -covered by a candidate representative pattern  $P_l$ , which is further  $\delta_g$ -covered by a final representative pattern  $P_g$ . An obvious result from Theorem 1 shows if we choose  $\delta_l + \delta_g = \delta$ , then  $P$  is guaranteed to be covered by  $P_g$ . Here we exploit a better assignment. We have:

$$D(P, P_l) = 1 - \frac{|T(P_l)|}{|T(P)|} \leq \delta_l$$

$$D(P_l, P_g) = 1 - \frac{|T(P_g)|}{|T(P_l)|} \leq \delta_g$$

The constraint given by the problem is:

$$D(P, P_g) = 1 - \frac{|T(P_g)|}{|T(P)|} \leq \delta$$

To achieve more compression, we would like the values of  $\delta_l$  and  $\delta_g$  to be as large as possible. This implies:

$$(1 - \delta_l) \times (1 - \delta_g) = 1 - \delta$$

We use  $\lambda$  to control the tradeoff between  $\delta_l$  and  $\delta_g$ , such that  $\delta_l = \lambda\delta$  and  $\delta_g = \frac{(1-\lambda)\delta}{1-\lambda\delta}$ .

We further discuss the selection of *min\_sup* for RPlocal and RPglobal:  $M_l$  and  $M_g$ . Since the original compression problem has parameter  $M$  and  $\delta$ , we have  $\hat{M} = (1 - \delta) \times M$  for representative patterns. The RPlocal step needs to keep the same *min\_sup* for representative patterns. Thus,

$$M_l = \frac{\hat{M}_l}{1 - \delta_l} = \frac{\hat{M}}{1 - \delta_l} = \frac{(1 - \delta)M}{1 - \lambda\delta}$$

All the frequent patterns (w.r.t.  $M$ ) are  $\delta_l$ -covered in the RPlocal step. To ensure that they are  $\delta$ -covered finally, the RPglobal step needs to  $\delta_g$ -cover all the patterns that possibly  $\delta_l$ -cover the frequent patterns (w.r.t.  $M$ ). Thus,

$$M_g = (1 - \delta_l) \times M = (1 - \lambda\delta) \times M$$

In conclusion, RPcombine takes three parameters:  $M, \delta, \lambda$ . The RPlocal step uses parameters  $M_l, \delta_l$  on the database, and the RPglobal step uses parameters  $M_g, \delta_g$  on the outputs of the first step.

## 5 Performance Study

To compare the proposed algorithms, a comprehensive performance study is conducted by testing our implementations on the data sets in the *frequent itemset mining dataset repository* [9]. All the algorithms were implemented in C++, and all the experiments were conducted on an Intel Pentium-4 2.6GHz system with 1GB RAM. The system ran Linux with the 2.6.1 kernel and gcc 3.3.2.

We summarize the methods to be compared as follows. In the FPClose method, we generate all the closed frequent patterns w.r.t.  $M$  (we use FPClose package[11], which is the winner of FIMI workshop 2003 [9]). In the RPglobal method, we first use FPClose to get all the closed frequent itemsets with *min\_sup*  $\hat{M} = M \times (1 - \delta)$ , then use RPglobal to find a set of representative patterns covering all the patterns with *min\_sup*  $M$ . In the RPlocal method, we directly compute all the representative patterns from database. For the clear presentation, we temporarily exclude the RPcombine method in the first three groups of experiments. One can imagine that the performance of RPcombine fits in somewhere between RPlocal and RPglobal. We will return to RPcombine later in the section.

### 5.1 Number of Presentative Patterns

The first set of experiments compare three algorithms w.r.t. the number of output patterns. We select *accidents*, *chess*, *connect* and *pumsb\_star* data sets [9].

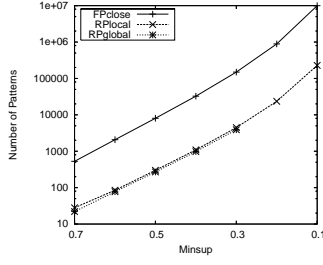


Figure 10: Number of Output Patterns w.r.t.  $min\_sup$ , Accidents Data Set

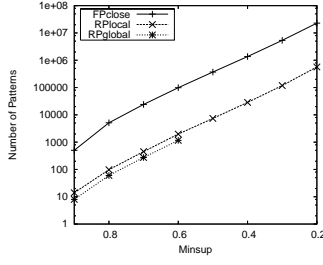


Figure 11: Number of Output Patterns w.r.t.  $min\_sup$ , Chess Data Set

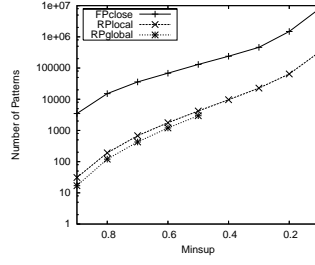


Figure 12: Number of Output Patterns w.r.t.  $min\_sup$ , Connect Data Set

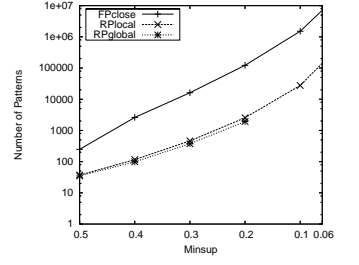


Figure 13: Number of Output Patterns w.r.t.  $min\_sup$ , Pumsb\_star Data Set

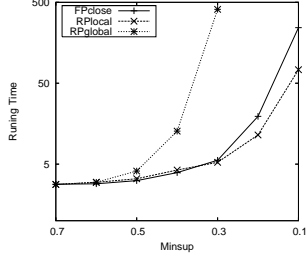


Figure 14: Running Time w.r.t.  $min\_sup$ , Accidents Data Set

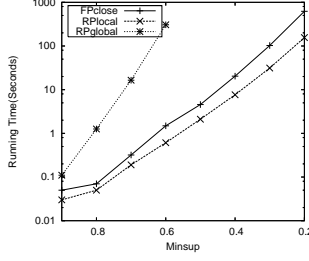


Figure 15: Running Time w.r.t.  $min\_sup$ , Chess Data Set

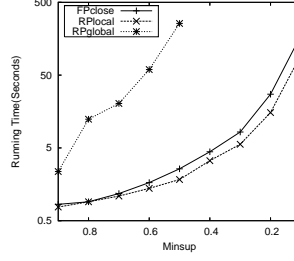


Figure 16: Running Time w.r.t.  $min\_sup$ , Connect Data Set

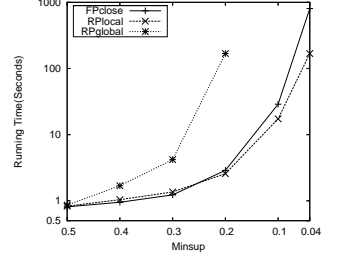


Figure 17: Running Time w.r.t.  $min\_sup$ , Pumsb\_star Data Set

For each data, we vary the value of  $min\_sup$  as the percentage of the number of total transactions and fix  $\delta = 0.1$  (we think it is a reasonably good compression quality).

The results are shown from Fig. 10 to Fig. 13. We have the following observations: First, both RPglobal and RPlocal are able to find a subset of representative patterns, which is almost two orders of magnitude less than the whole collection of the closed patterns; Second, although RPlocal outputs more patterns than RPglobal, the performance of RPlocal is very close to RPglobal. Almost all the outputs of RPlocal are within two times of RPglobal. The results of RPglobal are partial in that when minimum support becomes low, the number of closed patterns grows very fast, the running times of RPglobal exceed the time limit (30 minutes).

## 5.2 Running Time

The corresponding running time of the three methods are shown from Fig. 14 to Fig. 17. The times for RPglobal include FPClose procedure.

The results show that RPglobal does not scale well w.r.t. the number of patterns, and is much slower than RPlocal. Comparing FPClose and RPlocal, we observe that although RPlocal examines more patterns than FPClose (i.e., RPlocal examines the patterns with  $min\_sup \hat{M}$ , while FPClose only examines the patterns with  $min\_sup M$ ), RPlocal runs faster than FPClose, especially when  $min\_sup$  is low.

We further investigate the benefit of *closed pruning*. Fig. 18 shows the results on *pumsb\_star* data set, with

three configurations: FPClose, RPlocal ( $\delta = 0.1$ ) with and without *closed pruning*. We observe that even without *closed pruning*, RPlocal is more efficient than FPClose. This is because in RPlocal, the number of representative patterns is much less than the number of closed patterns. As a result, both the construction and query on RP-trees are more efficient. We use two integers as *closed\_index* and the improvement by applying *closed pruning* is significant. When  $M = 0.04$ , the *closed pruning* version runs three times faster than the version without *closed pruning*, and four times faster than FPClose. At that time, the number of frequent items is 173.

## 5.3 Distribution of Representative Patterns

The distributions of representative patterns w.r.t. pattern lengths and pattern supports are shown from Fig. 19 to Fig. 22. We use *accidents* data set, with  $min\_sup = 0.4$ . In order to get a high-level summary of support distributions, we group supports into 10 buckets. The bucket id is computed by  $\lfloor \frac{10 \times support}{\#transactions} \rfloor$ .

Fig. 19 shows the distributions w.r.t. the pattern lengths for three methods: FPClose, RPglobal and RPlocal ( $\delta = 0.1$ ). We observe that the overall shape of the distributions of RPglobal and RPlocal are similar to the shape of closed patterns. RPglobal and RPlocal have certain shifts to longer length because the nature of the compression problem favors larger itemsets. Fig. 20 compares the distributions (w.r.t. pattern supports) of close itemsets, maximal itemsets, and representative itemsets by RPglobal and RPlocal ( $\delta = 0.2$ ). While the

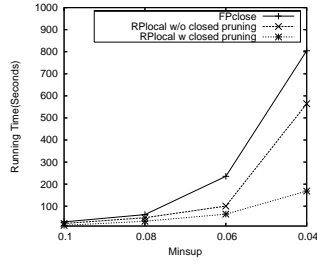


Figure 18: Running Time w.r.t closed pruning, Pumsb\_star Data Set

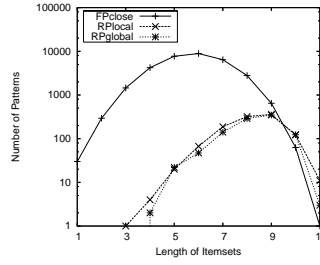


Figure 19: Distribution of Patterns w.r.t. Pattern Length, Accident Data Set

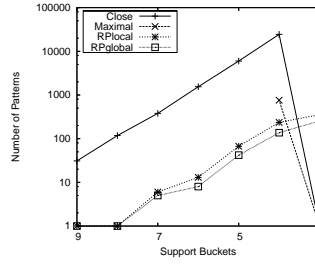


Figure 20: Distribution of Patterns w.r.t. Support Buckets, Accident Data Set

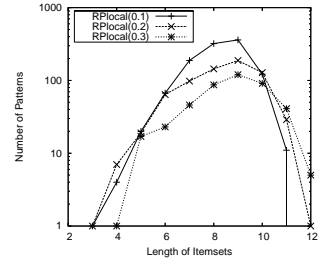


Figure 21: Distribution of Patterns w.r.t. Pattern Length, Accident Data Set

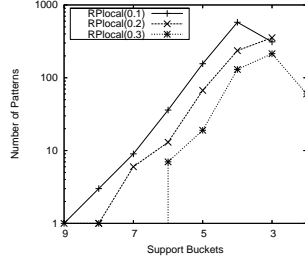


Figure 22: Number of Output Patterns w.r.t. Support Buckets, Accidents Data Set

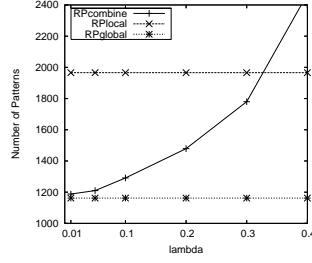


Figure 23: Number of Output Patterns w.r.t.  $\lambda$ , Chess Data Set

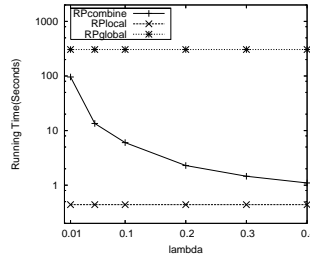


Figure 24: Running Time w.r.t.  $\lambda$ , Chess Data Set

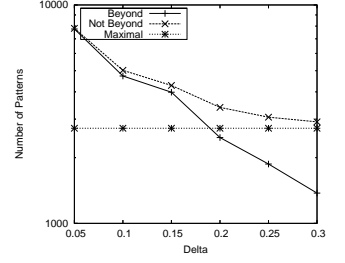


Figure 25: Number of Output Patterns, Accidents Data Set

maximal itemsets catch the boundary supports only, RPlglobal and RPlclose are able to get a reasonable distribution which is similar to the original closed itemsets. These suggest that both RPlglobal and RPlclose achieve high quality compressions.

We also run RPlclose with different  $\delta$  from 0.1 to 0.3. Fig. 21 and Fig. 22 show the pattern distributions w.r.t. lengths and supports. As we expected, the number of representative patterns decreases when the value of  $\delta$  increases, because a larger value of  $\delta$  enables a representative pattern to cover more patterns. Increasing the value of  $\delta$  also shifts the distributions of the patterns to longer and lower support patterns.

#### 5.4 Additional Tests

We examine the performance of RPcombine w.r.t. the different values of  $\lambda$ . Fig. 23 shows the final number of representative patterns by RPlglobal, RPlclose and RPcombine on *chess* data set (with  $M = 0.6, \delta = 0.1$ ). Fig. 24 shows the corresponding running time. The times of RPcombine are the sums of local and global steps. The  $\lambda$  for RPcombine is varied from 0.01 to 0.4. When  $\lambda$  is small (*i.e.*, 0.01), the local step reports more candidates and the global step takes more time, but the compression quality is better. The compression quality degrades as  $\lambda$  increases, the number of representative patterns is even larger than RPlclose when  $\lambda = 0.4$ . This is because at that time,  $M_g$  decreases a lot in order to guarantee all the original frequent patterns are covered, and  $\delta_g$  also decreases. As

a result, the global step needs to cover more pattern with tighter quality measure. In most applications, we suggest to choose  $\lambda = 0.1$ .

The final experiment is designed to verify the benefit to allow the support of representative patterns to be beyond of  $\min\_sup M$ . We compare the number of output patterns with three different options: (1) the *Beyond* case where the  $\min\_sup$  of representative patterns is  $\hat{M}$ ; (2) the *Not Beyond* case where the  $\min\_sup$  of representative patterns is  $M$ ; and (3) maximal patterns with  $\min\_sup M$ . We use the *accident* data set, varying  $\delta$  from 0.05 to 0.3, while fixing  $M$  as 0.3. The results in Fig. 25 show that the *beyond* case gets fewer number of representative patterns, especially in the case when  $\delta$  is large, while the *not beyond* case has maximal patterns as its lower bound.

## 6 Discussion

In this section, we discuss several related issues. First, to approximate a collection of frequent patterns, people always favor the more succinct compression. However, the explosive output pattern size restricts the application of most advanced algorithms. The RPlclose method can be used as sampling procedure as we did in RPcombine, since it is efficient and achieves considerable compression. Second, the compressed pattern sets generated by our method can be used for queries of finding approximate supports. We can construct an RP-tree with all the representative patterns. The query process is similar to the *coverage check*-

ing, except that in *coverage checking*, the query pattern comes with its support, while here, the support is unknown. Among all the patterns which subsume the query pattern, we report the maximum support  $C$ . The support of the query pattern is bounded in  $[C, \frac{C}{1-\delta}]$ .

## 7 Conclusions

We have considered the problem of compressing frequent patterns. The problem was shown to be NP-Hard. Several methods have been proposed. The RP-global method has theoretical bound, and works well on small collections of frequent patterns. The RP-local method is quite efficient, and preserves reasonable compression quality. We also discuss a combined approach, RPcombine, to balance the quality and efficiency.

There are several interesting directions that we are considering for future work: (1) Compression of sequence, graph and structure patterns; (2) Using compressed patterns for associations, correlations and classifications; and (3) Compressing frequent patterns over incrementally updated data (e.g., data streams).

## 8 Appendix: Proof of Theorem 4

**Proof.** For the simplicity of presentation, we prove the theorem in the set-covering framework. Let the sequence of probe elements be  $e_1, e_2, \dots, e_l$  ( $l = |C_l|$ ), the sets selected by the local method be  $S_1, S_2, \dots, S_l$ . Similar to the approaches in [7], we assign a cost 1 to each set which is selected by the local method, distribute this cost evenly over the elements covered for the first time. If  $e$  is covered for the first time by  $S_i$ , then  $c_e = \frac{1}{|S_i - (S_1 \cup S_2 \cup \dots \cup S_{i-1})|}$ , and  $|C_l| = \sum c_e$ . The cost assigned to the optimal cover is  $\sum_{S \in C^*} \sum_{e \in S} c_e$ . Since each element is in at least one set in the optimal cover, we have

$$|C_l| = \sum c_e \leq \sum_{S \in C^*} \sum_{e \in S} c_e = \sum_{S \in C^*} K_S$$

Where  $K_S = \sum_{e \in S} c_e = \sum_{S' \in C_l} \frac{|R(S') \cap S|}{|R(S')|}$  and  $R(S_i) = S_i - (S_1 \cup S_2 \cup \dots \cup S_{i-1})$ .

Let the minimum sets covering all the probe elements be  $\{M_1, M_2, \dots, M_T\}$ . We further assume that set  $M_i$ ,  $i \in \{1, 2, \dots, T\}$ , covers probe elements  $E_i = \{e_{i_1}, e_{i_2}, \dots, e_{i_p}\}$  ( $i_1 < i_2 < \dots < i_p$ ), where  $E_i \cap E_j = \emptyset$ , ( $\forall i \neq j$ ), and  $\cup_{i=1}^T E_i = \{e_1, e_2, \dots, e_l\}$ .

Let the set (selected by the local method) associated with probe element  $e_{i_1}$  be  $S_{i_1}$ . Since  $S_{i_1}$  is one of the current largest sets which cover probe pattern  $e_{i_1}$ , we have  $|R(S_{i_1})| \geq |R(M_i)|$ . Because  $e_{i_1}, e_{i_2}, \dots, e_{i_p}$  are in the order of probe elements, these elements must have not been covered at the time when  $e_{i_1}$  is selected as probe element. Thus,  $|R(M_i)| \geq p$ , and we conclude  $|R(S_{i_1})| \geq p$ . Similarly, we have  $|R(S_{i_2})| \geq p - 1$ ,  $|R(S_{i_3})| \geq p - 2, \dots, |R(S_{i_p})| \geq 1$ .

For all  $S' \in C_l$ , assume  $S'_1, S'_2, \dots, S'_l$  is in in ascending order of  $|R(S')|$ .  $K_S$  achieves the maximum value if we distribute the elements of  $S$  first in set  $S'_1$  fully, then  $S'_2, \dots$ , until  $S'_{k+1}$ . Since distributing fully on  $S'$  means  $\frac{|R(S') \cap S|}{|R(S')|} = 1$ , we have  $K_S \leq k + 1$ .

Evenly distribute the first  $k$   $S'$  into  $T$  buckets, and assign minimum  $|R(S')|$  value for them, we have,

$$\begin{aligned} |S| &\geq |R(S'_1)| + |R(S'_2)| + \dots + |R(S'_k)| \\ &\geq T \times \sum_{i=1}^k i = T \times \frac{\frac{k}{T} \times (\frac{k}{T} + 1)}{2} > \frac{k^2}{2T} \end{aligned}$$

We have  $K_S \leq k + 1 < \sqrt{2T|S|} + 1$ , thus

$$\begin{aligned} |C_l| &\leq |C^*| \times (\max_{S \in C^*} (K_S)) \\ &< |C^*| \times (\sqrt{2T \times \max_S |S|} + 1) \end{aligned}$$

## References

- [1] F. Afrati, et al. Approximating a Collection of Frequent Sets. KDD'04.
- [2] R. Agrawal, et al. Mining association rules between sets of items in large databases. SIGMOD'93.
- [3] R. Agrawal and R. Srikant. Fast algorithm for mining association rules. VLDB'94.
- [4] R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95.
- [5] R. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98.
- [6] S. Brin, et al. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97.
- [7] T. Cormen, et al. Introduction to Algorithms, Second Edition. MIT Press. 2001.
- [8] G. Dong and J. Li. Efficient mining of emerging patterns: Discovering trends and differences. KDD'99.
- [9] Frequent Itemset Mining Dataset Repository. <http://fimi.cs.helsinki.fi/data/>
- [10] K. Gouda and M. Zaki. Efficiently Mining Maximal Frequent Itemsets. ICDM'01.
- [11] G. Grahne and J. Zhu. Efficiently Using Prefix-trees in Mining Frequent Itemsets. IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03).
- [12] J. Han, et al. Efficient mining of partial periodic patterns in time series database. ICDE'99.
- [13] J. Han, et al. Mining Frequent Patterns without candidate generation. SIGMOD'00.
- [14] B. Lent, et al. Clustering association rules. ICDE'97.
- [15] H. Mannila, et al. Discovery of frequent episodes in event sequences. Data Mining and Knowledge Discovery, 1:259-289, 1997
- [16] N. Pasquier, et al. Discovering Frequent Closed Itemsets for Association Rules. ICDT'99.
- [17] C. Silverstein, et al. Scalable techniques for mining causal structures. VLDB'98.
- [18] J. Wang, et al. Closet+: Searching for the Best Strategies for Mining Frequent Closed Itemsets. KDD'03.
- [19] M. Zaki and C. Hsiao. Charm: An Efficient Algorithm for Closed Itemset Mining. SDM'02.