# Efficient Algorithms for Mining Frequent Itemsets with Constraint

4 authors:

Anh Tran
University of Dalat

**11** PUBLICATIONS   **65** CITATIONS

SEE PROFILE

Hai Duong
University of Dalat, Vietnam

**9** PUBLICATIONS   **72** CITATIONS

SEE PROFILE

Tin C Truong
University of Dalat

**21** PUBLICATIONS   **117** CITATIONS

SEE PROFILE

Bac Le
Ho Chi Minh City University of Science

**62** PUBLICATIONS   **300** CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:

Project    Sequence Mining View project

Project    Vision-based traffic surveillance View project

# Efficient Algorithms for Mining Frequent Itemsets with Constraint

Anh N. Tran
University of Dalat,
Dalat, Vietnam
anhtn@dlu.edu.vn

Hai V. Duong
University of Dalat,
Dalat, Vietnam
haidv@dlu.edu.vn

Tin C. Truong
University of Dalat,
Dalat, Vietnam
tintc@dlu.edu.vn

Bac H. Le
University of Natural Science
Ho Chi Minh, Vietnam
lhbac@fit.hcmus.edu.vn

*Abstract*—**An important problem of interactive data mining is "to find frequent itemsets contained in a subset $C$ of set of all items on a given database". Reducing the database on $C$ or incorporating it into an algorithm for mining frequent itemsets (such as Charm-L, Eclat) and resolving the problem are very time consuming, especially when $C$ is often changed. In this paper, we propose an efficient approach for mining them as follows. Firstly, it is necessary to mine only one time from database the class $LG_A$ containing the closed itemsets together their generators. After that, when $C$ is changed, the class of all frequent closed itemsets and their generators on $C$ is determined quickly from $LG_A$ by our algorithm MINE_CG_CONS. We obtain the algorithm MINE_FS_CONS to mine and classify efficiently all frequent itemsets with constraint from that class. Theoretical results and the experiments proved the efficiency of our approach.**

*Closed itemsets, frequent itemsets, constraint, generators, eliminable itemsets*

## I. INTRODUCTION

Currently, Internet makes the real changes in the ways human thinks and does. People access to Internet to get useful information from it. Normally, the data of websites for users is obtained and is saved in the tables (or databases). The number of attributes (items) is often enormous. However, for a while, they only take care of a set of attributes (called the constraint). To show immediately to the users the knowledge mined from them such as the frequent itemsets or association rules is very important. In [3, 5, 8, 9, 13, 14] some authors researched on mining frequent itemsets and association rules from the standpoint of the user's interaction with the system. *They studied mining frequent itemsets with many different kinds of constraints*. Nguyen et al. [9] proposed an architect including domain, class and SQL-style aggregate constraints. Some categories of constraints such as anti-monotone, monotone, and succinct have been integrated into the mining process [9]. In [13], Pei et al. proposed the concept of convertible constraints and considered pushing them into the mining progress of the FP-growth algorithm. Srikant et al. [14] considered the problem of integrating constraints that are Boolean expressions over the presence or absence of items in the association rules. Bayardo et al. [3] restricted the problem of mining association rules in two constraints of the consequent and the

minimum improvement. In [5], Cong and Liu proposed an technique based on the concept of tree boundary to utilize previous mining results for reducing the mining time. They considered tightening and relaxing constraints such as increasing and decreasing supports.

This paper concentrates on solving the problem of to find frequent itemsets contained in a subset $C$ of the set of all items from a given database (called the problem of mining frequent itemsets with constraint $C$). A simple approach is to reduce the database on $C$ and after that to mine them by an algorithm used widely such as Apriori [1], *Eclat* [20], *Charm-L* [18], FP-growth [8], etc. This approach is not efficient because the constraints are often changed. A different one is to filter the output of those algorithms in a post-processing step to determine frequent itemsets with constraint. It is also not efficient because their outputs are usually enormous. In [14], Srikant et al. showed that we should incorporate $C$ into the mining process. They modified the apriori candidate generation procedure to only count candidates that contain selected items. However, there are still many candidates generated. Furthermore, when $C$ is changed, users must run the algorithm. It is very time consuming. In [8], Han et al. also suggested incorporating $C$ into the mining FP-Tree. However, they did not propose the algorithm to do it. We also incorporate constraint $C$ into Charm-L and Eclat (well-known algorithms for mining frequent closed itemsets and all frequent itemsets) to mine frequent itemsets with constraint. This approach will be compared to our approach.

Recently, we [15, 16] showed the structure of each class of equivalent frequent itemsets having the same closure based on the generators and eliminable subsets of that closure. Based on it, we approach to the problem of mining frequent itemsets with constraint as follows. Firstly, it is necessary to mine only one time from the database the class $LG_A$ containing the closed itemsets and their generators. After that, when $C$ is changed, the class $FLG_C$ of frequent closed itemsets and their generators is determined quickly from $LG_A$. Using a unique representation of frequent itemset, we derive completely, directly, non-repeatedly from $FLG_C$ all frequent itemsets with constraint. The mining script is shown as follows:

1) To mine only one time the class $LG_A$

*2) User selects the constraint $C$ and the minimum support threshold s:*

*(2.1) to determine quickly the class $FLG_A^s$ of the frequent closed itemsets (and their generators) with respect to s from $LG_A$ in the first mining time; otherwise, from $FLG_A^{s_{max}}$ (was saved before), where $s_{max}$ is maximum such that $s_{max} \leq s$,*

*(2.2) from $FLG_A^s$ , our algorithm MINE_CG_CONS (based on the propositions 2 and 3) is used to exploit directly the class $FLG_C$,*

*(2.3) from $FLG_C$ the class $FS_C$ of all frequent itemsets with constraint $C$ is mined quickly by our algorithm MINE_FS_CONS (using theorem 4).*
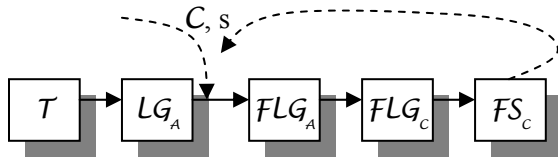
*3) Return the step 2.*



Figure 1.   Mining frequent itemsets with constraint.

The method is suitable when $C$ is usually changed. Indeed, the size of the class of all frequent closed itemsets and their generators is much smaller than the one of all frequent itemsets, especially on the dense databases. With the small values of minimum support thresold, this class can be still mined and saved in main memory by Charm-L and MinimalGenerators [18]. The response of the system in early times is often slow because the size of $LG_A$ is big. After the first period, the classes of $FLG_A^{s_1}$, $FLG_A^{s_2}$, .., $FLG_A^{s_n}$ are saved in the system (the values of $s_j$, j=1, .., n are distributed regularly on [0, 1], $0 \leq s_1 < s_2 < .. < s_n \leq 1$). Thus, we only need to select the frequent closed itemsets and generators that those supports exceed threshold s directly on $FLG_A^{s_{max}}$, where $s_{max} = \max \{s_j \mid s_j \leq s, j=1, .., n\}$. For the threshold s given by user, the corresponding threshold $s_{max}$ is usually closed to it. Therefore, the time to exploit $FLG_A^s$ is often small. Using some simple operators on the itemsets in $FLG_A^s$, we can mine directly the class $FLG_C$ of the frequent closed itemsets and generators restricted on $C$. The class $FS_C$ is partitioned into the disjoint equivalence classes. Each class contains frequent itemsets on $C$ that have the same closure. So, they can be mined concurrently by parallel algorithms. A unique representation of frequent itemsets based on frequent closed itemset (represented to that class) and its generators is indicated to derive directly, non-repeatedly all frequent itemsets on $C$ from $FLG_C$.

The rest of the paper is organized as follows. Section 2 recalls some primitive concepts and results. This section also proposes a unique representation of itemsets. The main results are in section 3. In it, we obtain the algorithm for determining quickly all frequent closed itemsets and their generators with constraint $C$. The efficient algorithm to mine all frequent itemsets from them is also figured out. Sections 4 and 5 contain the experimental results and conclusions.

## II.   PRIMITIVE CONCEPTS AND RESULTS

### A.   Primitive concepts

Given set $\mathcal{O}$ contained records or transactions of a database $\mathcal{T}$ and $\mathcal{A}$ contained attributes or items related to each of transaction $o \in \mathcal{O}$ and $\mathcal{R}$ is a binary relation in $\mathcal{O}x\mathcal{A}$. Consider two operators: $\lambda:2^{\mathcal{O}} \rightarrow 2^{\mathcal{A}}$, $\rho:2^{\mathcal{A}} \rightarrow 2^{\mathcal{O}}$ determined as $(\rho(\varnothing) := \mathcal{O},\ \lambda(\varnothing) := \mathcal{A})$:

$$\lambda(O) = \{a \in \mathcal{A} \mid (o, a) \in \mathcal{R},\ \forall o \in O\},\ \forall O \subseteq \mathcal{O},$$
$$\rho(A) = \{o \in \mathcal{O} \mid (o, a) \in \mathcal{R},\ \forall a \in A\},\ \forall A \subseteq \mathcal{A}.$$

Defining closure operator h in $2^{\mathcal{A}}$ [4] by: $h = \lambda \circ \rho$, we say that h(A) is the closure of itemset $A \subseteq \mathcal{A}$. If A = h(A), A is called *closed itemset*. The class of all closed itemsets is denoted as $CS$. The *support* of itemset A is defined as the probability of the ocurrence of a transaction containing A on $\mathcal{O}$: $s(A) = |\rho(A)|/|\mathcal{O}|$. Denoted that $s_0$ is minimum support, $s_0 \in [1/|\mathcal{O}|; 1]$, if $s(A) \geq s_0$ then A is called *frequent itemset* [1]. Let $FS$, $FCS$ be the classes of all frequent itemsets and all frequent closed itemsets.

For two non-empty itemsets G, A: $\varnothing \neq G \subseteq A \subseteq \mathcal{A}$, G is called a *generator* of A [11] iff: $h(G)=h(A)$ and $(\forall \varnothing \neq G' \subset G \Rightarrow h(G') \subset h(G))$. Let $G(A)$ be the class of all generators of A. Let $LG_A$ and $FLG_A$ [1] be the classes of all closed itemsets together their generators on $\mathcal{A}$ and all elements in $LG_A$ that are frequent with respect to $s_0$. In $2^{\mathcal{A}}$, an itemset R is called eliminable in S [15, 16] iff $R \subset S$ and $\rho(S) = \rho(S \backslash R)$. Let $N(S)$ denote the class of all *eliminable itemsets* in S, $N^*(S) := N(S) \backslash \{\varnothing\}$, we have [15]: $N(S) = \{A: A \subseteq S \backslash G, G \in G(S)\}$.

TABLE I.        DATABASE 1

| Trans ID | Items |
|----------|-------|
| 1 | aceg |
| 2 | acfh |
| 3 | adfh |
| 4 | bceg |

*Example 1.* Let us consider database 1 in Table I, with minimum support $s_0 = \frac{1}{4}$, used in all next examples of this paper. From the definitions of $\lambda$ and $\rho$, we have: $\lambda(\{1, 4\}) = ceg$, $\rho(ceg) = \{1, 4\}$ and then, h(ceg)=ceg. So ceg is a frequent closed itemset with the support $|\rho(ceg)|/|\mathcal{O}| = \frac{1}{2}$.

---

[1] For briefly, we write $FLG_A^s$ simply $FLG_A$ .

This itemset contains two generators e, g because h(e)=h(g)=h(ceg)=ceg.

### B. Structure of Itemsets

In this part, the class of all itemsets is partitioned into the disjoint equivalence classes. The elements of an equivalence class have the same closure and can be derived from that closure and its generators. Thus, we only need to mine the class (with the small size) of frequent closed itemsets (and generators). When it is necessary, we can derive to users the frequent itemsets in a class that they are interested in.

*Definition 1* [15] *(Equivalence relation $\sim_h$ over the class of all itemsets $2^A$)*: $\forall$ A, B $\in$ $2^A$:

$$A \sim_h B \Leftrightarrow h(A) = h(B).$$

*Theorem 1* [15] *(A partition of $2^A$)*: Relation $\sim_h$ partitions $2^A$ into the disjoint equivalence classes. Each class contains itemsets that have the same closure. The equivalence class containing A is denoted as [A].

$$2^A = \sum_{A \in CS} [A] \text{ and } FS = \sum_{A \in FCS} [A].$$

Based on this partition, we can exploit independently each equivalence class. The elements in a class have the same support so we only compute and save it once.

*Theorem 2* [15] *(Representation of itemset)*: For every itemset A such that $\varnothing \neq A \in CS$:

$$X \in [A] \Leftrightarrow \exists G_0 \in G(A), \exists X' \in N(A): X = G_0 + X'. \text{ }^2$$

Denoted $N(S, G) := \{A: A \subseteq S \backslash G, G \in G(S)\}$, it is obvious that: $N(S) = \bigcup_{G \in G(S)} M(S, G)$. For $G_1, G_2 \in G(S)$, $G_1 \neq G_2$, the intersection of $N(S, G_1)$ and $N(S, G_2)$ can be not empty. The above representation of a itemset in a class can be not unique because the representation of an eliminable itemset is not unique.

*Example 2.* Let us consider equivalence class [X], where X=aceg, $G(X) = \{ae, ag\}$. We have: $N^*(X, ae) = \{cg, c, g\}$, $N^*(X, ag)=\{ce, c, e\}$, $N^*(X, ae) \cap N^*(X, ag) = \{c\} \neq \varnothing$ and $N^*(X) = N^*(X, ae) \cup N^*(X, ag) = \{cg, c, g, e, ce\}$. Then, from theorem 2, [X] = {ae, aeg, aec, aecg, ag, agc}, in which, aeg can be represented by two ways: aeg = ae+g = ag+e.

### C. Unique Representation of Itemset by Generator and Eliminable Itemset

The process of deriving of all itemsets of an equivalence class using theorem 2 can make the duplication because the representation of an itemset is not unique. Theorem 3 shows a unique representation of itemset, in other words, based on it, all itemsets in the same class can be derived non-repeatedly (as a result, quickly), completely.

$\forall X \in CS$, let us call $X_U = \bigcup_{X_i \in G(X)} X_i$ , $X_{U,i}=X_U \backslash X_i$, $X_\_ =X \backslash X_U$, $IS(X) := \{X'=X_i+X'_i+X^\sim \mid X_i \in G(X), X^\sim \subseteq X_\_, X'_i \subseteq X_{U,i}, i=1 \text{ or } (i>1: X_k \not\subset X_i+X'_i, \forall k: 1 \leq k<i)\}$.

*Theorem 3 (Unique representation of itemset by generator and eliminable itemset)*: We have:

a. [X] = IS(X).

b. All itemsets of IS(X) are derived non-repeatedly.

*Proof*:

(a) .*"$\subseteq$"*: If $X' \in [X]$, assume that i is the minimum index such that $X_i \in G(X)$, $X''_i \subseteq X \backslash X_i$ and $X' = X_i+X''_i$. Let $X'_i = X''_i \cap X_U$, $X^\sim = X''_i \backslash X_U$, then $X'_i \subseteq X_{U,i}$, $X^\sim = X' \backslash X_U \subseteq X_\_$ and $X' = X_i+X'_i+X^\sim$. Assume that there exists the index k such that $1 \leq k<i$, $X_k \in G(X)$, $X_k \subseteq X_i+X'_i$ then $X'=X_k+X''_k$, where $X''_k=X'_k+X^\sim$ and $X'_k=(X_i+X'_i) \backslash X_k \subseteq X \backslash X_k$, $X^\sim \subseteq X \backslash X_k$. Thus, $X''_k \subseteq X \backslash X_k$: it is absurd!

.*"$\supseteq$"*: If $X' \in IS(X)$, there exists $X_i \in G(X)$, $X^\sim \subseteq X_\_ \subseteq X \backslash X_i$, $X'_i \subseteq X_{U,i} \subseteq X \backslash X_i$: $X'=X_i+X'_i+X^\sim$. Let $X''= X'_i+X^\sim \in N(X)$, then $X'=X_i+X''$, thus $X' \in [X]$ by theorem 2.

(b) Assume that there exists i, k such that $i>k \geq 1$ and $X_i+X'_i+X^\sim_i \equiv X_k+X'_k+X^\sim_k$, where: $X_i$, $X_k \in G(X)$; $X^\sim_i$, $X^\sim_k \subseteq X_\_$; $X'_i \subseteq X_{U,i}$, $X'_k \subseteq X_{U,k}$. Since $X_k \cap X^\sim_i=\varnothing$, so $X_k \subset X_i+X'_i$ (the equality do not occur because $X_i$ and $X_k$ are two different generators of X). It contradicts to the selection of index i! $\square$

*Example 3.* Let us consider class [X] where X=aceg, $G(X) = \{X_1=ae, X_2=ag\}$, we have: $X_U=aeg$, $X_{U,1}=g$, $X_{U,2}=e$, $X_\_=c$. By theorem 3, itemset $X'=aceg \in IS(X)$ is generated uniquely as follows: $X'=X_1+X'_1+X^\sim$ where $X'_1=g \subseteq X_{U,1}$, $X^\sim=c \subseteq X_\_$. By theorem 2, X' has two duplicate representations: X'=ae+cg=ag+ce. If the condition "i>1: $X_k \not\subset X_i+X'_i, \forall k: 1 \leq k<i$" is absent, then duplicate X' is generated once again: $X'=X_2+X'_2+X^\sim$, where $X'_2=e \subseteq X_{U,2}$ and $X_1 \subset X_2+X'_2$. Similarly, all itemsets of [X] = IS(X) = {ae, aeg, aegc, aec, ag, agc} are derived non-repeatedly.

From theorem 3, the algorithm **GEN_ITEMSETS** is obtained (see Fig. 2) to generate *non-repeatedly all itemsets* in *each equivalence class* [X], X $\in$ CS.

### III. MINING FREQUENT ITEMSETS WITH CONSTRAINT

As the discussion in introduction, to mine frequent itemsets on $C$ with minimum support threshold $s_0$, firstly, without the general, we need to determine the class $FLG_A$ of all frequent (with respect to $s_0$) closed itemsets and their generators from $LG_A$. After that, it is quickly to mine the class $FLG_C$ of all frequent closed itemsets and generators restricted on $C$ from $FLG_A$. That bases on some relations between closed itemsets and generators of $FLG_C$ and the corresponding ones of $FLG_A$. Finally, the partition of the

---

$^2$ The symbol + is denoted as the union of two disjoint sets.

class of all frequent itemsets with constraint $C$ allows us to use the algorithm *GEN_ITEMSETS* for deriving quickly them.

---

$<IS(X), s(X)>$  **GEN_ITEMSETS** $(X, s(X), G(X))$:

1. $IS(X) = \varnothing$; $X_U = \bigcup_{X_i \in G(X)} X_i$ ; $X\_ = X \backslash X_U$;

2. **for each** $(i=1; X_i \in G(X); i{++})$ **do** {
3.     $X_{U,i} = X_U \backslash X_i$;
4.     **for each** $(X'_i \subseteq X_{U,i})$ do {
5.       IsDuplicate = false;
6.       **for** $(k=1; k<i; k{++})$
7.         **if** $(X_k \subset X_i + X'_i)$ **then** {
8.           *IsDuplicate* = true; **break**;
9.         }
10.       **if** *(not(IsDuplicate))* **then**
11.         **for each** $(X^\sim \subseteq X\_)$ do
12.           $IS(X) = IS(X) + \{X_i + X'_i + X^\sim\}$;
13.     }
14. }
15. **return** $<IS(X), s(X)>$;

---

Figure 2.   GEN_ITEMSETS, the algorithm to generate non-repeatedly all itemsets in class [X].

### A. Mining frequent closed itemsets and their generators with constraint

We will define again the operators $\lambda$, $\rho$ and $h$ over constraint $C$ and figure out the relation between them with the corresponding ones over $A$. From that, the algorithm for mining quickly $FLG_C$ from $LG_A$ is indicated.

*Definition 2 (The Galois connection operators over constraint C)*: For every $C \in 2^A \backslash \{\varnothing\}$, let us consider operators: $\rho_C : 2^C \to 2^O$, $\lambda_C : 2^O \to 2^C$ and $h_C : 2^C \to 2^C$ defined as follows: $\forall \varnothing \neq C' \subseteq C$, $\varnothing \neq O \subseteq O$,

$$\rho_C(C') = \{o \in O: (o, a) \in R, \forall a \in C'\}, \rho_C(\varnothing) := O,$$
$$\lambda_C(O) = \{a \in C: (o, a) \in R, \forall o \in O\}, \lambda_C(\varnothing) := C,$$
$$h_C = \lambda_C \circ \rho_C.$$

An itemset $C' \subseteq C$ is called closed itemset on $C$ iff $h_C(C') = C'$. The class of all frequent itemsets on $C$ is denoted as $FS_C$. The class $FCS(C)$ contains all frequent closed itemsets on $C$.

*Proposition 1*: For every $C \in 2^A \backslash \{\varnothing\}$, $\varnothing \neq C' \subseteq C$, $O \subseteq O$, the following statements are true:

  a. $\rho_C(C') = \rho(C')$, so $s(C') = |\rho(C')| = |\rho_C(C')|$,
  b. $\lambda_C(O) = \lambda(O) \cap C$,
  c. $h_C(C') = h(C') \cap C$.

*Proof:* Obviously from definitions of $\rho$, $\rho_C$, $\lambda$, $\lambda_C$, $h$ and $h_C$. $\square$

Proposition 1.c enables us to determine frequent closed itemset on $C$ by intersecting $C$ with each frequent closed itemset (on $A$) of $FLG_A$. This way can make the duplications. In other words, a frequent closed itemset on $C$ can be derived many times.

*Example 4.* Let us consider database 1, we have: $FLG_A$={acfh, aceg, adfh, bceg, *afh*, ceg, *ac*, c, *a*}. Then, by proposition 1.c, with $C$=abde, $FLG_C$={ae, ad, be, e, a}. Some of its elements are derived many times, for example: a = $acfh \cap C = afh \cap C = ac \cap C = a \cap C$.

Based on a condition over generators, proposition 2 is obtained to eliminate the duplication in generating frequent closed itemsets restricted on $C$.

*Proposition 2 (Generating non-repeatedly all frequent closed itemsets with constraint C):* Let us call $FCS_C := \{C'=L \cap C \mid L \in FCS, \exists L_i \in G(L): L_i \subseteq C\}$, we have:

  a. $FCS_C = FCS(C)$.
  b. *All elements of $FCS_C$ are generated non-repeatedly.*

*Proof:* a. "$\subseteq$": $\forall \varnothing \neq C' \in FCS_C$: $C' = L \cap C \subseteq C$, $h(C') \subseteq h(L) = L$. Then, $C' \subseteq h_C(C') = h(C') \cap C \subseteq L \cap C = C'$. Therefore $C' = h_C(C')$, i.e., $C' \in FCS(C)$.

"$\supseteq$": $\forall \varnothing \neq C' \in FCS(C)$: $C' = h_C(C') = h(C') \cap C = L \cap C$, where $L := h(C') \in FCS$. Let $C_i \in G(C')$ (there always exists), then: $h(C_i)=h(C')=L=h(L)$ and $\forall C_0 \subset C_i$ then $h(C_0) \subset h(C_i) = L$. Thus, $C_i \in G(L)$ and $C_i \subseteq C' \subseteq C$, i.e., $C' \in FCS_C$. We conclude that $C' \in FCS_C$.

b. Assume that, with k=1, 2, $C'_k = L_k \cap C \in FCS_C$, let $L_k \in FCS$, $L_{k,0} \in G(L_k)$, $L_{k,0} \subseteq C$ such that $C'_1 \equiv C'_2$ and $L_1 \neq L_2$. We have $L_{k,0} \subseteq L_k \cap C$, $L_k = h(L_{k,0}) \subseteq h(L_k \cap C) \subseteq h(L_k) = L_k$. Therefore, $h(C'_k) = L_k$, $\forall k=1, 2$ and $L_1 = L_2$: it is absurd! $\square$

In the next step, we will show how to determine the generators of frequent closed itemsets on $C$.

*Definition 3 (The generators of C' restricted on C)*: For every $G, C'$: $\varnothing \neq G \subseteq C' \subseteq C$, $G$ is called a generator of $C'$ on $C$ iff:

$$h_C(G)=h_C(C') \text{ and } (\forall \varnothing \neq G' \subset G \Rightarrow h_C(G') \subset h_C(G)).$$

The set of all generators of $C'$ on $C$ is denoted as $G_C(C')$.

*Proposition 3 (Determining the generators with constraint C):* $\forall C'=L \cap C \in FCS_C$:

$$G_C(C') = G(C') = \{L_i \in G(L): L_i \subseteq C'\}.$$

*Proof:* It is obvious because $\rho_C(C') = \rho(C')$. $\square$

From propositions 2 and 3, the algorithm **MINE_CG_CONS** is indicated to *mine quickly the class* $FLG_C := \{<C', s(C'), G_C(C')> \mid C' \in FCS_C\}$ from $LG_A$.

$\mathcal{FLG}_C$  **MINE_CG_CONS** ($LG_A$, $C$, $s_0$):
1. $\mathcal{FLG}_C = \varnothing$;
2. **for each** (<L, s(L), $\mathcal{G}$(L)> $\in LG_A$) **do**
3.   **if** (s(L) $\geq s_0$) **then**   // L $\in \mathcal{FLG}_A$
4.     **if** ($\exists$ L$_i \in \mathcal{G}$(L) **and** $C \supseteq$ L$_i$) **then** {
        // not to generate repeatedly
5.       C'   = L $\cap$ $C$;
6.       $\mathcal{G}_C$(C') = {L$_i \in \mathcal{G}$(L) | L$_i \subseteq$ C'};
7.       $\mathcal{FLG}_C = \mathcal{FLG}_C$ + <C', s(L), $\mathcal{G}_C$(C')>;
8.     }
9. **return** $\mathcal{FLG}_C$;

Figure 3.   *MINE_CG_CONS*, the algorithm to generate non-repeatedly all frequent closed itemsets and their generators with constraint $C$.

*Example 5.* The process of mining frequent closed itemsets and generators on $C=abde$ from $LG_A=${<L, s(L), $\mathcal{G}$(L)>} is shown in Table II (C'=L$\cap C$).

TABLE II.       AN EXAMPLE OF MINING FREQUENT CLOSED ITEMSETS AND THEIR GENERATORS WITH CONSTRAINT

| L | L$_i \in \mathcal{G}$(L) | s(L) | $C \supseteq$L$_i$ | C' | $\mathcal{G}_C$(C') | s(C') |
|---|---|---|---|---|---|---|
| acfh | cf, ch | ¼ | | | | |
| aceg | ae, *ag* | ¼ | ae | ae | ae | ¼ |
| adfh | d | ¼ | d | ad | d | ¼ |
| bceg | b | ¼ | b | be | b | ¼ |
| afh | f, h | ½ | | | | |
| ceg | e, *g* | ½ | e | e | e | ½ |
| ac | ac | ½ | | | | |
| c | c | ¾ | | | | |
| a | a | ¾ | a | a | a | ¾ |

### B. Mining all frequent itemsets with constraint

Here, we will partition the class $\mathcal{FS}_C$ of all frequent itemsets restricted on $C$ into the disjoint equivalence classes. Each class contains the itemsets having the same closure with the frequent closed itemset represented to that class. Thus, it is correct to use the efficient algorithm *GEN_ITEMSETS* for mining quickly $\mathcal{FS}_C$ from $\mathcal{FLG}_C$.

*Definition 4 (Equivalence relation over $2^C$)*: $\forall$ A, B $\in 2^C$:
$$A \sim_C B \Leftrightarrow h_C(A) = h_C(B).$$

*Theorem 4 (Partition and representation of $\mathcal{FS}_C$)*: The equivalence relation $\sim_C$ partitions $\mathcal{FS}_C$ into the disjoint equivalence classes. Each class contains the frequent itemsets having the same closure:
$$\mathcal{FS} = \sum_{C' \in \mathcal{FCS}(C)} [C']_{\sim_C} = \sum_{C' \in \mathcal{FCS}_C} \mathcal{IS}(C').$$

*Proof:* This theorem is consequence of theorems 1, 3 and the proposition 2.                                    □

The algorithm **MINE_FS_CONS** mines and classifies quickly all frequent itemsets on $C$ from $LG_A$.

$\mathcal{FS}_C$   **MINE_FS_CONS** ($LG_A$, $C$, $s_0$):
1. $\mathcal{FS}_C = \varnothing$;
2. $\mathcal{FLG}_C =$ **MINE_CG_CONS** ($LG_A$, $C$, $s_0$);
3. **for each** (<C', s(C'), $\mathcal{G}_C$(C')> $\in \mathcal{FLG}_C$) **do** {
4.   <$\mathcal{IS}$(C'), s(C')> =
        **GEN_ITEMSETS** (C', s(C'), $\mathcal{G}_C$(C'));
5.   $\mathcal{FS}_C = \mathcal{FS}_C$ + {<$\mathcal{IS}$(C'), s(C')>};
        // classify $\mathcal{FS}_C$
6. }
7. **return** $\mathcal{FS}_C$;

Figure 4.   *MINE_FS_CONS*, the algorithm to generate non-repeatedly all frequent itemsets on $C$.

*Example 6.* The processes of mining from $LG_A = ${X=<L, s(L), $\mathcal{G}$(L)>} all frequent itemsets restricted on $C_1$=abde and $C_2$=abceg are figured out in Tables III and IV, where $\mathcal{FLG}_C = ${Y = <C', s(C'), $\mathcal{G}_C$(C')>}.

TABLE III.       MINING ALL FREQUENT ITEMSETS ON $C_1$

| X $\in LG_A$ | Y $\in \mathcal{FLG}_C$ | <$\mathcal{IS}$(C'), s(C')> |
|---|---|---|
| acfh, ¼, {cf, ch} | | |
| aceg, ¼, {ae, ag} | ae, ¼, {ae} | {ae}, ¼ |
| adfh, ¼, {d} | ad, ¼, {d} | {d, da}, ¼ |
| bceg, ¼, {b} | be, ¼, {b} | {b, be}, ¼ |
| afh, ½, {f, h} | | |
| ceg, ½, {e, g} | e, ½, {e} | {e}, ½ |
| ac, ½, {ac} | | |
| c, ¾, {c} | | |
| a, ¾, {a} | a, ¾, {a} | a, ¾ |

TABLE IV.       MINING ALL FREQUENT ITEMSETS ON $C_2$

| X $\in LG_A$ | Y $\in \mathcal{FLG}_C$ | <$\mathcal{IS}$(C'), s(C') > |
|---|---|---|
| acfh, ¼, {cf, ch} | | |
| aceg, ¼, {ae, ag} | ac*eg*, ¼, {ae, ag} | {ae, aec, aeg, aegc, ag, agc}, ¼ |
| adfh, ¼, {d} | | |
| bceg, ¼, {b} | b*ceg*, ¼, {b} | {b, bc, be, bg, bce, bcg, beg, bceg}, ¼ |
| afh, ½, {f, h} | | |
| ceg, ½, {e, g} | c*eg*, ½, {e, g} | {e, ec, eg, egc, g, gc}, ¼ |
| ac, ½, {ac} | ac, ½, {ac} | {ac}, ½ |
| c, ¾, {c} | c, ¾, {c} | {c}, ¾ |
| a, ¾, {a} | a, ¾, {a} | {a}, ¾ |

### IV. EXPERIMENTAL RESULTS

The following experiments were performed on a 2.93 GHz Pentium(R) Dual-Core CPU E6500 with 1.94GB of

RAM, running Linux, Cygwin. Algorithms were coded in $C^{++}$. The code of Zaki [22] is used to run Charm-L, MinimalGenerators and Eclat. Four databases in [21] are used during these experiments. They have been used as benchmark for testing mining algorithms. Table V shows their characteristics.

TABLE V.        DATABASE CHARACTERISTICS

| Database (DB) | # Records | # Items | Average size |
|---|---|---|---|
| Mushroom (M) | 8124 | 119 | 23 |
| Connect (Co) | 67557 | 129 | 43 |
| Chess (Ch) | 3196 | 75 | 37 |
| Pumsb (P) | 49046 | 7117 | 50 |

As the discussion in the introduction, Srikant et al. incorporated $C$ into the mining process by modified the apriori candidate generation procedure. In this experiment, to compare with our approach in mining frequent itemsets with constraint from $LG_{A}$, we incorporate $C$ into the Charm-L and Eclat (well-known algorithms for mining frequent itemsets). This is done easily by choosing only frequent items that are in $C$ to work in the next steps of those algorithms. Those new versions are called $C$–Charm, and $C$–Eclat.

The items of the constraints are selected from the set $A^F$ of all frequent (corresponding to the minimum support MS) items of $A$ with the ratios of ¼, ½ and ¾. We have the constraints with the sizes of $l_1=¼*|A^F|$, $l_2=½*|A^F|$ and $l_3=¾*|A^F|$. For each $l_i$, $C$ is constructed from two subsets: $C = C_1+C_2$. ==In the reality, the constraints that users are interested in usually contain the high-support items. Thus, we will sort all items by the order of their supports.== To determine a constraint $C$ with the size $l_i$, firstly, we construct the first subset $C_1$ of $C$ containing $[p*l_i]$ items randomly selected from the set of high-support items in $A^F$, where $p \in [0; 1]$; the remained part $C_2$ of $C$ contains $[(1-p)*l_i]$ randomly selected items from $A^F \backslash C_1$. For experiments at here, we set $p = 0.5$ and consider two constraints for each $l_i$.
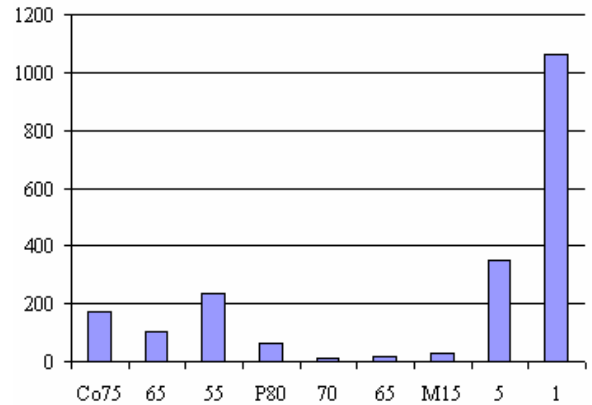
We will do two comparisons. Firstly, in Table VI, we compare the average time for mining frequent closed itemsets and their generators with constraint $C$ by our algorithm *MINE_CG_CONS* (shown in column $T_C^O$) to the one by $C$–CharmGen (column $T_C^C$) upon P, M, Co, and Ch. $C$–CharmGen is the combination of $C$–Charm (for mining frequent closed itemsets with constraint) and MinimalGenerators [19] (for determining their generators). The reduction in the mining time is shown in column $RT_C^O$ ($RT_C^O = T_C^C/T_C^O$). With the different minimum support thresholds, we see that it is drastic, ranging from a factor of *60* to **316** times! That reduction plays an important role in mining quickly association rule with constraints because, as

the discussion on [15, 6], all association rules can be mined quickly from frequent closed itemsets and their generators.

TABLE VI.        MINING FREQUENT CLOSED ITEMSETS AND GENERATORS WITH CONSTRAINT: *MINE_CG_CONS* VS *C-CHARMGEN*.

| DB | MS (%) | $T_C^C$ | $T_C^O$ | $RT_C^O$ |
|---|---|---|---|---|
| Co | 75 | 1.07 | 0.005 | *213* |
|  | 65 | 1.29 | 0.013 | *99* |
|  | 55 | 3.49 | 0.026 | *134* |
| Ch | 65 | 1.21 | 0.011 | *115* |
|  | 55 | 6.08 | 0.055 | *111* |
|  | 50 | 10.93 | 0.111 | *98* |
| P | 80 | 1.50 | 0.010 | *145* |
|  | 70 | 13.74 | 0.115 | *120* |
|  | 65 | 39.22 | 0.243 | *161* |
| M | 15 | 0.16 | 0.003 | *60* |
|  | 5 | 0.84 | 0.003 | ***316*** |
|  | 1 | 4.80 | 0.018 | *264* |

In the second, the time for mining all frequent closed itemsets restricted on $C$ by our algorithm *MINE_FS_CONS* is compared to the one by $C$–Eclat. We did experiments on on three databases (that have many items) Co, P, M. The reduction in the mining time by our approach is drastic. It ranges from a factor of 14 to *1063* times and is shown in Fig. 5.



Figure 5.    The reduction in mining time all frequent itemsets with constraint: *MINE_FS_CONS* vs *C–Eclat*.

For database P, the reductions are small because the average size of transactions is small compared with the number of all items. Then, practically (users are usually interested in high-support items), we should consider the constraints containing many high-support items. Indeed, the bigger of the number of high-support items (corresponding with the big values of p) are, the bigger reductions become. Table VII showed that. In addition, the output of our algorithm *MINE_FS_CONS*  is classified into disjoint classes. All frequent itemsets with constraint in a class have

24

the same closure and support. Thus, when it is necessary, we only access without the need to compute them.

TABLE VII.  THE REDUCTIONS IN MINING TIME FREQUENT ITEMSETS WITH CONSTRAINT: *DATABASE P*

| MS / p | 80 | 70 | 65 |
|---|---|---|---|
| ½ | 61 | 13 | 15 |
| ¾ | 51 | 26 | 36 |
| 4/5 | 70 | 37 | 65 |

## V.  CONCLUSIONS

This paper proposed an approach to mine and classify efficiently frequent itemsets with constraint $C$ on a given database, especially, when $C$ is often changed. The correctness and efficiency of the approach were ensured by the theoretical results. The corresponding algorithms were obtained and were tested on benchmark databases. In future, based on this approach, we will research on the problem of mining association rules with constraint.

## REFERENCES

[1]  R. Agrawal, and R. Srikant, "Fast algorithms for mining association rules", Proceeding of the 20th International Conference on Very Large Data Bases, 1994, pp. 478-499.

[2]  R.J. Bayardo, "Efficiently Mining Long Patterns from Databases," Proceedings of the SIGMOD Conference, 1998, pp. 85–93.

[3]  R.J. Bayardo, R. Agrawal, and D. Gunopulos, "Constraint-Based Rule Mining in Large, Dense Databases," Data Mining and Knowledge Discovery, Kluwer Academic Publishers, vol. 4, No. 2/3, 2000, pp. 217–240.

[4]  H.T. Bao, "An approach to concept formation based on formal concept analysis," IEICE Trans. Infor. and systems, E78-D, No. 5, 1995.

[5]  G. Cong, and B. Liu, "Speed-up Iterative Frequent Itemset Mining with Constraint Changes," ICDM, 2002, pp. 107-114.

[6]  B. Ganter, and R. Wille, Formal Concept Analysis: Mathematical Foundations, Springer-Verlag, 1999.

[7]  J. Han, J. Pei, and J. Yin, "Mining frequent itemsets without candidate generation," Proceedings of SIGMOID'00, 2000, pp. 1-12.

[8]  J. Han, J. Pei, Y. Yin, and R. Mao, "Mining frequent patterns without candidate generation: a frequent-pattern tree approach, " Data mining and knowledge discovery, no 8, 2004, pp. 53-87.

[9]  R.T. Nguyen, V.S. Lakshmanan, J. Han, and A. Pang, "Exploratory Mining and Pruning Optimizations of Constrained Association Rules," Proceedings of the 1998 ACM-SIG-MOD Int'l Conf. on the Management of Data, pp. 13-24.

[10]  N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, "Efficient mining of association rules using closed item set lattices," Information systems, vol. 24, no. 1, 1999, pp. 25-46.

[11]  N. Pasquier, R. Taouil, Y. Bastide, G. Stumme, and L. Lakhal, "Generating a condensed representation for association rules," J. of Intelligent Information Systems, vol. 24, no. 1, 2005, pp. 29-60.

[12]  J. Pei, J. Han, and R. Mao, "CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets," Proceedings of the DMKDWorkshop on Research Issues in Data Mining and Knowledge Discovery, 2000, pp. 21–30.

[13]  J. Pei, J. Han, and V.S. Lakshmanan, "Pushing Convertible Constraints in Frequent Itemset Mining," Data Mining and Knowledge Discovery, no 8, 2004, pp. 227–252.

[14]  R. Srikant, Q. Vu, and R. Agrawal, "Mining association rules with item constraints," Proceeding KDD'97, pp. 67-73.

[15]  T.C. Tin, and T.N. Anh, "Structure of set of association rules based on concept lattice," SCI 283, Advances in Intelligent Information and Database Systems, Springer-Verlag, Berlin Heidelberg, 2010, pp. 217-227.

[16]  T.C. Tin, T.N. Anh, and T. Thong, "Structure of Association Rule Set based on Min-Min Basic Rules," Proceedings of the 2010 IEEE-RIVF International Conference on Computing and Communication Technologies, pp. 83-88.

[17]  R. Wille, "Concept lattices and conceptual knowledge systems," Computers and Math. with App., no. 23, 1992, pp. 493-515.

[18]  M.J. Zaki, and C.J. Hsiao, "Efficient algorithms for mining closed itemsets and their lattice structure," IEEE Trans. Knowledge and data engineering, vol. 17, no. 4, 2005, pp. 462-478.

[19]  M.J. Zaki, "Mining non-redundant association rules," Data mining and knowledge discovery, no. 9, 2004, pp. 223-248.

[20]  M.J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li, "New algorithms for fast discovery of association rules," Proc. 3rd Int. Conf. on Knowledge Discovery and Data Mining (KDD'97), pp. 283–296.

[21]  Frequent Itemset Mining Dataset Repository (FIMDR), http://fimi.cs.helsinki.fi/data/, acessed 2009.

[22]  http://www.cs.rpi.edu/~zaki/www-new/pmwiki.php/Software/Software#patutils, acessed 2010.