

Fast and Memory Efficient Mining of Frequent Closed Itemsets

Claudio Lucchese
Salvatore Orlando
Raffaele Perego

From TKDE06

Outline

- Introduction
- Memory-Efficient Duplicate Detection and Pruning
- DCI_Closed Algorithm
- Performance Analysis
- Conclusion

Introduction

- **Dense data:**
Contain strongly correlated items and long frequent patterns
- Such data sets are, in fact, very hard to mine, while the number of frequent itemsets grows up very quickly as the minimum support threshold is decreased.

Introduction(cont.)

- **Closed Itemsets**
Given an itemset $T \subseteq D$, and $I \subseteq I$ and we define
 $f(T) = \{ i \in I \mid \forall t \in T, i \in t \}$
 $g(S) = \{ t \in D \mid \forall i \in I, i \in t \}$
- **An itemset I is said to be closed if and only if**

$$c(I) = f(g(I)) = f \circ g(I) = I$$

Introduction(cont.)

- $\text{min_supp} = 1,$

TID	Items
1	B D
2	A B C D
3	A C D
4	C

(a) (b)

Introduction(cont.)

- **Browsing the search space**
 Lemma 1. Given two itemsets X and Y ,if $X \subset Y$ and $\text{supp}(X) = \text{supp}(Y)$, then $c(X) = c(Y)$
- Therefore, given a generator X, if we find an already mined closed itemsets Y that set-include X, where the supports of Y and X are identical, we can conclude that $c(X) = c(Y)$. In this case, we also say that Y subsumes X. If this holds, we can safely prune the generator X without computing its closure. Otherwise, we have to compute $c(X)$ in order to obtain a new closed itemset.

Introduction(cont.)

- We could in fact mine all the closed itemsets by computing the closure of just this single representative itemset for each equivalence class, without generating any duplicate. Let us call representative itemsets *closure generators*.
- Other algorithms use a different technique, which we call *closure climbing*.
- For example, the closed itemset $\{A,B,C,D\}$ of the figure could be mined twice since it can be obtained as the closure of two minima elements of its equivalence class, namely, $\{A,B\}$ and $\{B,C\}$.

Introduction(cont.)

- Given an itemset X and an item $i \in I$, $g(X) \subseteq g(i)$
 $\Leftrightarrow i \in c(X)$
- From the above lemma, we have that if $g(X) \subseteq g(i)$, then $i \in c(X)$. Therefore, by performing this inclusion check for all the items i not included in X , we can incrementally compute $c(X)$.

Memory-Efficient Duplicate Detection and Pruning(cont.)

- For example, the closed itemsets $\{A,C,D\}$ has four such generators, namely, $\{A\}$, $\{A,C\}$, $\{A,D\}$, and $\{C,D\}$.
- Denote with symbol $<$ the usual lexicographic total order between two ordered itemsets, in turn, defined on the basis of R .

Memory-Efficient Duplicate Detection and Pruning(cont.)

- A generator of the form $X=Y \cup i$, where Y is a closed itemset and $i \notin Y$, is said to be order-preserving iff either $c(X) = X$ or $i < c(X) \setminus X$.
- Example of Figure, we have that $\{A\} = \psi \cup \{A\}$ is an order-preserving generator of the closed itemset $\{A,C,D\}$, while $\{C,D\} = \{C\} \cup \{D\}$ is not an order-preserving generator for the same closed itemset.

Memory-Efficient Duplicate Detection and Pruning(cont.)

- In order to mine all the closed itemsets by avoiding redundancies, we compute the closure of order-preserving generators only and prune the others.
- **Theorem 1.**
 For each closed itemset $\bar{Y} \neq c(\psi)$, there exists a sequence of n items $i_0 < i_1 < \dots < i_{n-1}$, $n \geq 1$, such that $\langle \text{gen}_0, \text{gen}_1, \dots, \text{gen}_{n-1} \rangle = \langle Y_0 \cup i_0, Y_1 \cup i_1, \dots, Y_{n-1} \cup i_{n-1} \rangle$, where the various gen_i are order-preserving generators, with $Y_j \neq c(\psi)$, $j \in [0, n-1]$, $Y_{j+1} = c(Y_j \cup i_j)$, and $Y_n = c(\bar{Y})$.

Memory-Efficient Duplicate Detection and Pruning(cont.)

- **Corollary 1.**
 For each closed itemset $\bar{Y} \neq c(\psi)$, the sequence of order-preserving generators of Theorem 1 is unique.
- **Example:**
 For the closed itemset $\bar{Y} = \{A,B,C,D\}$, we have $Y_0 = c(\psi) = \psi$, $\text{gen}_0 = \psi \cup \{A\}$, $Y_1 = c(\text{gen}_0) = \{A,C,D\}$, $\text{gen}_1 = \{A,C,D\} \cup \{B\}$, and, finally, $Y = c(\text{gen}_1)$.

Memory-Efficient Duplicate Detection and Pruning(cont.)

□ **Detecting Order-Preserving Generator**

□ **Definition 3.**

Given a generator $gen = YUi$, where Y is a closed itemset and $i \notin Y$, we define $pre\text{-set}(gen)$ as follows:

$$pre\text{-set}(gen) = \{ j \mid j \in I, j \notin gen, \text{ and } j < i \}.$$

□ **Lemma 3.**

Let $gen = YUi$, i be a generator where Y is a closed itemset and $i \notin Y$. If $\exists j \in pre\text{-set}(gen)$ such that $g\{gen\} \subset g\{j\}$, then gen is not order-preserving.

DCI_CLOSED Algorithm

□ DCI_CLOSED starts by scanning the input data set D to determine the frequent single items $F1 \subseteq I$ and builds the bitwise vertical data set VD containing the various tidlists $g(i)$.

□ After this first step, DCI_CLOSED decides whether VD corresponds to either a dense or a sparse data set. Since VD is bitwise, if the percentage of 1s is large, the data set is soon classified as dense.

DCI_CLOSED Algorithm(cont.)

Algorithm 1 DCI_CLOSED pseudocode

```

1: procedure DCI_CLOSEDd(CLOSED.SET, PRE.SET, POST.SET)
2: while POST.SET ≠ ∅ do
3:   i ← mini(POST.SET)
4:   POST.SET ← POST.SET \ i
5:   new_gen ← CLOSED.SET ∪ i           ▷ Build a new generator
6:   if supp(new_gen) ≥ min_supp and
7:     ¬is_dup(new_gen, PRE.SET) then
8:     CLOSED.SETnew ← new_gen
9:     POST.SETnew ← ∅
10:    for all j ∈ POST.SET do           ▷ Compute closure of new_gen
11:      if g(new_gen) ⊆ g(j) then
12:        CLOSED.SETnew ← CLOSED.SETnew ∪ j
13:      else
14:        POST.SETnew ← POST.SETnew ∪ j
15:    end for
16:    Write Out CLOSED.SETnew and its support
17:    DCI_CLOSEDd(CLOSED.SETnew, PRE.SET, POST.SETnew)
18:  end if
19: end while
20: end procedure

```

DCI_CLOSED Algorithm(cont.)

```

22: function is_dup(new_gen, PRE.SET)           ▷ Duplicate check
23: for all j ∈ PRE.SET do
24:   if g(new_gen) ⊆ g(j) then
25:     return TRUE                             ▷ new_gen is not order preserving
26:   end if
27: end for
28: return FALSE
29: end function

```

DCI_CLOSED Algorithm(cont.)

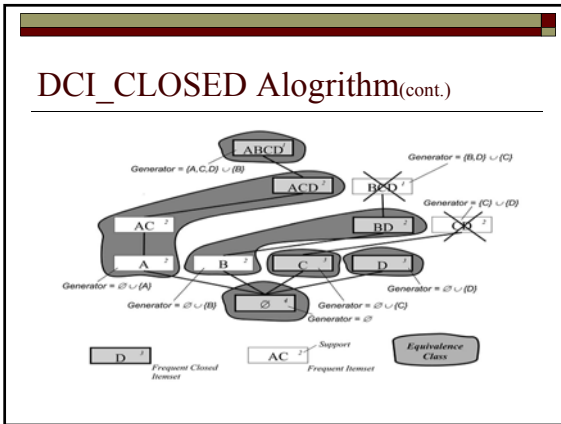
□ Once $c(\psi)=\psi$ is found, four generators can be constructed by adding a single item to $c(\psi)$, namely, $\{A\}$, $\{B\}$, $\{C\}$, and $\{D\}$. Suppose we first compute the closure of $gen=\psi \cup \{A\}=\{A\}$. Note that, since no items precede A in the lexicographic order, then its PRE_SET is empty and, thus, we can conclude that gen is order-preserving. $DCI_CLOSED_d()$ checks if $g(A)$ is set-included in $g(j)$, $\forall j \in POST_SET$ (i.e., $g(B)$, $g(C)$, and $g(D)$), and discovers that $c(A)=\{A,C,D\}$.

□ $DCI_CLOSED_d()$ is then recursively called, with parameters $CLOSED_SET = \{A,C,D\}$, $POST_SET = \{B\}$, while PRE_SET is still empty. $CLOSED_SET = \{A,C,D\}$ is thus extended with B (its $POST_SET$), so obtaining a new generator $gen = \{A,C,D\} \cup \{B\} = \{A,B,C,D\}$. Since PRE_SET is empty, this generator is order-preserving by definition, but is also closed because $POST_SET$ is now empty.

DCI_CLOSED Algorithm(cont.)

□ After this first recursive exploration, $DCI_CLOSED_d()$ starts solving another independent subproblem by exploring generator $gen=\psi \cup \{B\}=\{B\}$, where $PRE_SET = \{A\}$ and $POST_SET = \{C,D\}$.

□ Finally, $DCI_CLOSED_d()$ starts exploring the last generator $gen = \psi \cup \{D\}=\{D\}$, where $PRE_SET = \{A,B,C\}$ and $POST_SET = \psi$. Since gen is order-preserving (this is checked by comparing $g(D)$ with $g(A)$, $g(B)$, and $g(C)$, i.e., with its PRE_SET), it is not pruned. But, we also can conclude that $\{D\}$ is also closed since $POST_SET = \psi$.

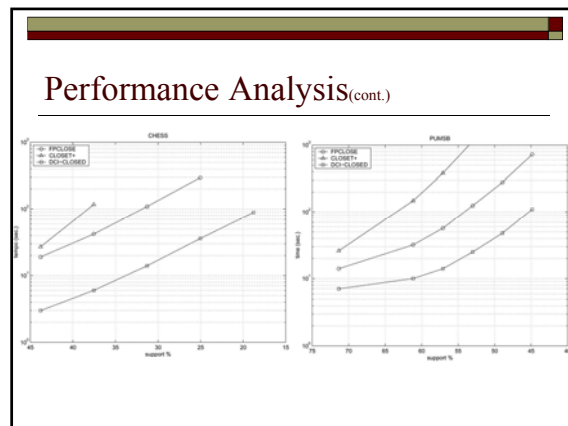
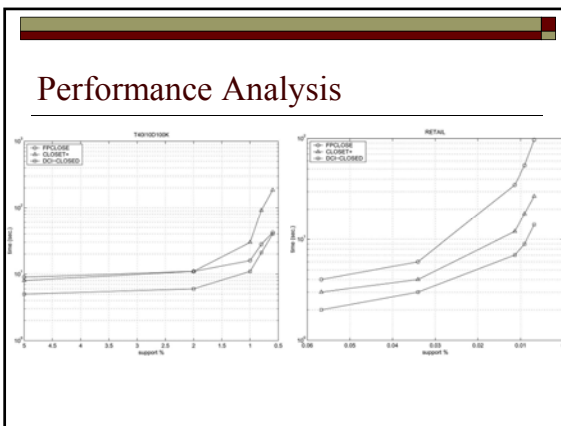


DCI_CLOSED Algorithm(cont.)

- Optimization Saving Bitwise Operations**
 - Data sets with highly correlated items

如果我們要mine的目標是x，那我們在每個columns與x無關的拿掉，但是這樣又會很花時間去check 每一個column，那我們就限制條件，限制只有跟Φ相關的itemset才可以執行。
 - Data sets with highly correlated items

此方發使用在mine dense data上，則是利用A Multi-Strategy Algorithm for Mining Frequent Sets以及 Adaptive and Resource-Aware Mining of Frequent Sets這兩篇paper的方式去作處理。



Conclusion

- In this paper, we have investigated the problem of efficiency in mining closed frequent itemsets from transactional data sets.
- Finally, it showed that allows dense data sets to also be effectively mined with the lowest possible support threshold