

6.1 Algorithm for determining frequent itemset and its support given a set of close frequent itemsets : Countlist = []

if $X = \emptyset$:

print ("not frequent itemset")

for each frequent closed itemset $f \in C$:

if $X \subset f$:

Countlist.append (support(f))

return Countlist.max() max(Countlist)

if Countlist == []

print ("not frequent itemset")

Note that: the support(f) corresponds to max in Countlist is the one with shortest length. The algorithm return the support of X if X is a frequent itemset, otherwise just print "not frequent itemset".

6.3

(a). Prove all non-empty subsets of a frequent itemset must be frequent. (Proof.)

Let A be a frequent itemset, min-sup be the minimum support. Let D be dataset, $|D|$ be number of all transactions, $A' \subseteq A$. any transaction

We have $\text{support-count}(A) \geq \text{min-Sup} \times |D|$
 $\text{contains } A \text{ also have } A'$

since $A' \subseteq A$ so $\text{support-count}(A') \geq \text{support-count}(A)$
 $\geq \text{min-Sup} \times |D|$

Hence, A' the subset of frequent itemset A must also be frequent.

(b). Since any transactions include s must also include s' where $s' \subset s$ but not the other way around. $\Rightarrow \text{Support-Count}(s') \geq \text{Support-Count}(s) \Rightarrow \text{Support}(s') \geq \text{Support}(s)$

This proves part (b).

(c). Since $\text{Confidence}(s' \Rightarrow l-s') = \frac{\text{Support}(l)}{\text{Support}(s')}$

Also $\text{confidence}(s \Rightarrow l-s) = \frac{\text{Support}(l)}{\text{Support}(s)}$

From part b, we know since $s' \subset s \Rightarrow \text{support}(s') \geq \text{support}(s)$

Hence. $\text{confidence}(s' \Rightarrow l-s') \leq \text{confidence}(s \Rightarrow l-s)$

(d). Since we know an itemset A is frequent in D . so $\text{Support-Count}(A) \geq \text{min-sup} \times |D|$, where min-sup is minimal support, $|D|$ is total # transactions of dataset D .

If we divide D into n non-overlapping dataset D_1, \dots, D_n

(Δ) Assume that A is not frequent in any of partition of D . $\Rightarrow \frac{\text{Support}(A \cap D_i)}{|D_i|} < \text{min-sup}$ for $i \in \{1, \dots, n\}$

$$\Rightarrow \text{Support}_{\{A\}}(D_i) \leq \min\text{-sup} \times |D_i|$$

for all $i \in \{1, \dots, n\}$

$$\Rightarrow \sum_{i=1}^n \text{Support}_{\{A\}}(D_i) = \sum_{i=1}^n \text{count}_{\{A\}}(D_i) = \text{Support}_{\{A\}}(D)$$

$$< \min\text{-sup} \times \sum_{i=1}^n |D_i| = \min\text{-sup} \times |D|$$

\therefore This contradicts (*) $\text{Support-Count}(A) \geq \min\text{-sup} \times |D|$ so our assumption (Δ) is wrong.

Hence, $\forall A \in C$ is frequent in at least one partition of D . (Proof by contradiction)

6.4. We actually only need to check $(k-2)$ subsets of length $(k-1)$ of C . Since we know C is generated from $(k-1)$ itemset L_1, L_2, \dots, L_k , there's no need to check frequent whether L_1, L_2 are in L_{k-1} . One way to do this is to pass L_1, L_2 as arguments to function has-infrequent-subset.

6.5. Since we know $\text{Conf}(S' \Rightarrow L-S') \leq \text{Conf}(S \Rightarrow L-S)$, if $S' \subset S$. We can improve association rule generation. If we already know a subset of S of L shows $\text{Conf}(S \Rightarrow L-S) < \min\text{-conf}$, then there's no need to check all the non-empty subsets of S . In this way, we don't need to check every subset of L whether each of them has $\text{conf} \geq \min\text{-sup}$.

Algorithm: Recursion

Input: L frequent itemset, minimal confidence (min-conf).

Output: itemsets in L with confidence \geq min-conf.

for each frequent itemset L
generating-itemset(L , L , min-conf)

Function generating-itemset (C : current subset
of L, L: original frequent itemset, min-conf).

K = length(C)

if ($K > 1$) then
Generate all $(K-1)$ subsets of L
for each $(K-1)$ subset m of L
if ($\text{support}(L) / \text{support}(m) \geq \text{min-conf}$)
print (" $m \Rightarrow (L - m)$ ")
generating-itemset-association(m, L),

This is much more efficient. we min-conf.
don't need to check every subset of L,
in fact, we only check those whose superset S
has ($S \Rightarrow L-S$)

6.6. min-sup = 60%, min-conf = 80%

(a). Part I: find frequent itemset using Apriori.

$\therefore \text{min-sup} = 60\% \therefore \text{min-support-count} = 3$.

$$L_1 = \{M, O, K, E, Y\}$$

$$C_2 = \{MO, MK, ME, MY, OK, OE, OY, KE, KY, EY\}$$

$$L_2 = \{MK, OK, OE, KE, KY\}$$

$$C_3 = \{OKE\}$$

$$L_3 = \{OKE\}$$

$$C_4 = \emptyset$$

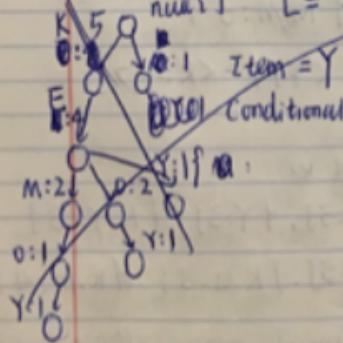
$$L_4 = \emptyset$$

so the set of frequent itemsets is:

$$\{M, O, K, E, Y, MK, OK, OE, KE, KY, OKE\}$$

Part II: find frequent itemset using FP-growth

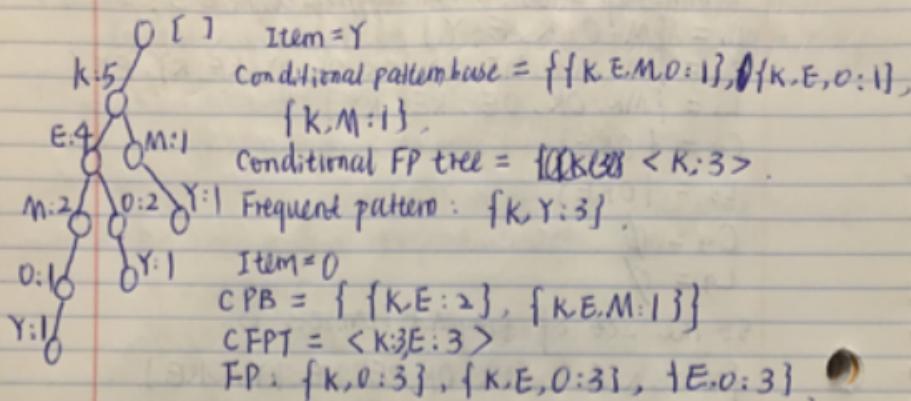
$$\text{L} = \{f(O:3), f(K:4), \{M:3\}, \{O:3\}, \{Y:3\}\}$$



$$\text{Conditional pattern base} = \{f(E, K, M, O:1), f(E, K, O:1)\}$$

Part II. find frequent itemset use FP-Growth.

$$L = \{\{K:5\}, \{E:4\}, \{M:3\}, \{O:3\}, \{Y:3\}\}$$



Hence, the frequent itemset is

$$\begin{aligned} & \{f\{K:5\}, f\{E:4\}, f\{M:3\}, f\{O:3\}, f\{Y:3\}, f\{K,Y:3\}, \\ & \{K,O:3\}, \{E,O:3\}, f\{K,E,O:3\}, f\{K,M:3\}, f\{K,E:4\}\} \end{aligned}$$

Summary: FP-growth is much more efficient since it only needs to scan database twice and we don't need to generate candidate sets frequently.

(b).

Let's take a look at the frequent itemset we found in part (a). Note: we only need to consider itemset with three items. Start from itemset with three items $\{K, E, O\}$: 3

We found. $\forall X \in E$ transaction \exists buys(X, D) \wedge

$$\textcircled{1}. \text{buys}(X, E) \Rightarrow \text{buys}(X, K) [60\%, 100\%]$$

$$\textcircled{2}. \text{buys}(X, D) \wedge \text{buys}(X, K) \Rightarrow \text{buys}(X, E) [60\%, 100\%]$$

~~for itemset with two items,~~ we found. $\forall X \in E$ transaction. $\text{buys}(X, K) \Rightarrow$

~~buys(X, E)~~ [80%],

$$\textcircled{3}. \text{buys}(X, E) \wedge \text{buys}(X, K) \Rightarrow \text{buys}(X, D) [60\%, 75\%]$$

Hence, $\textcircled{1}, \textcircled{2}$ are strong association matches requirement (form).

6.11.

When extending to Apriori and FP Growth by considering multiple occurrences of items, we need to treat ~~one~~ item with different counts as every different items. In our first scan of database for both methods, we will generate frequent single itemset by taking account to different counts. Then check if the minimal support is met. For Apriori, we can continue to

construct frequent 2-itemsets, 3-itemsets.
For FP-growth, when we project database,
we also have to project on itemsets where
each item can be associated with different
counts.