

# DATA MINING ASSIGNMENT - 5

Parampreet Singh

1(a) let  $s$  be a frequent itemset,

$\min\_sup$  be the minimum support

$D$  be the task relevant data, a set of database transactions

$|D|$  be the no. of transactions in  $D$

Since  $s$  is a frequent itemset  $support\_count(s) = \min\_sup \times |D|$

So,  $s'$  be any non empty set of  $s$   
Then any transaction containing itemset  $s$  will also contain itemset  $s'$ . Therefore  $support\_count(s') \geq support\_count(s) = \min\_sup \times |D|$

Thus,  $s'$  is also a frequent itemset

(b) let  $D$  be the task relevant data, a set of database transactions,  $|D|$  be the no. of transactions in  $D$ .

So,  $support(s) = \frac{support\_count(s)}{|D|}$   $s$  be the itemset

let  $s'$  be the non-empty subset of  $s$

$support(s') = \frac{support\_count(s')}{|D|}$



we know that  $\text{support}(s') \geq \text{support}(s)$

therefore the support of non-empty subset  $s'$  of itemset  $s$  will be greater than the support of  $s$ .

c) Given frequent itemset  $l$  and subset  $s$  of  $l$ ,

$$\text{confidence}(s \rightarrow (l-s)) = \frac{\text{support}(l)}{\text{support}(s)}$$

let  $s'$  be the any non-empty subset  $s$  then,

$$\text{confidence}(s' \rightarrow (l-s')) = \frac{\text{support}(l)}{\text{support}(s')}$$

we already know that  $\text{support}(s') \geq \text{support}(s)$

therefore,  $\text{confidence}(s' \rightarrow (l-s')) \leq \text{confidence}(s \rightarrow (l-s))$   
that is confidence of the rule  $s' \rightarrow l-s'$  cannot be more than the confidence of the rule  $s \rightarrow (l-s)$

Q. Given database has 5 transactions,  
 let  $\text{min\_sup} = 60\%$   
 $\text{min\_conf} = 80\%$

TID	Items-bought
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, V, C, K, Y}
T500	{I, O, O, K, I, E}

(a) Apriori Algorithm:-

So, No. of transactions = 5

$$\begin{aligned}\text{min support count} &= \frac{60}{100} \times (\text{total no. of transactions}) \\ &= \frac{60}{100} \times 5 = 3\end{aligned}$$

$$\begin{aligned}\text{min confidence count} &= \frac{80}{100} \times (\text{total no. of transactions}) \\ &= \frac{80}{100} \times 5 = 4\end{aligned}$$

$C_k$  = candidate itemset  
 $L_k$  = frequent itemset



we will generate the candidate itemset and also the frequent itemset by scanning the DB scans.

→ Need to remove the repeated item in same transaction based on association rule mining.

C<sub>1</sub>

- $\{m\} \rightarrow 1+1+1 = 3$
- $\{o\} \rightarrow 1+1+1 = 3$
- $\{N\} \rightarrow 1+1 = 2$
- $\{K\} \rightarrow 1+1+1+1+1 = 5$
- $\{E\} \rightarrow 1+1+1+1 = 4$
- $\{Y\} \rightarrow 1+1+1 = 3$
- $\{D\} \rightarrow 1 = 1$
- $\{A\} \rightarrow 1 = 1$
- $\{U\} \rightarrow 1 = 1$
- $\{C\} \rightarrow 1+1 = 2$
- $\{I\} \rightarrow 1 = 1$

L<sub>1</sub>

- $\{M\} \rightarrow 3$
- $\{O\} \rightarrow 3$
- $\{K\} \rightarrow 5$
- $\{E\} \rightarrow 4$
- $\{Y\} \rightarrow 3$

C<sub>2</sub>

- $\{M, O\} \rightarrow 1 = 1$
- $\{M, K\} \rightarrow 1+1 = 2$
- $\{M, E\} \rightarrow 1+1 = 2$
- $\{M, Y\} \rightarrow 1+1 = 2$
- $\{O, E\} \rightarrow 1+1+1 = 3$

L<sub>2</sub>

- $\{M, K\} - 3$
- $\{O, E\} - 3$
- $\{O, K\} - 3$
- $\{K, E\} - 4$
- $\{K, Y\} - 3$

$$\{o, k\} \rightarrow 3$$

$$\{o, y\} \rightarrow 2$$

$$\{k, E\} \rightarrow 4$$

$$\{k, y\} \rightarrow 3$$

$$\{E, y\} \rightarrow 2$$

$L_3$

$$\{m, k, o\} \Rightarrow 1 = 1$$

$$\{m, k, E\} \rightarrow 1+1 = 2$$

$$\{m, k, y\} \rightarrow 1+1 = 2$$

$$\{k, o, E\} \rightarrow 1+1+1 = 3$$

$$\{k, o, y\} \rightarrow 1+1 = 2$$

$$\{k, E, y\} \rightarrow 1+1 = 2$$

$$\{o, E, y\} \rightarrow 1+1 = 2$$

$L_3$

$$\{k, o, E\} \rightarrow 3$$

So the total frequent itemsets are

$$\{m, o, k, E, y, mk, oE, ok, kE, ky, koE\}$$

FP growth Algorithm:-

$$\text{minimum support} = \frac{60}{100} \times 5 = 3$$

Item sup

$$m = 3$$

$$o = 3$$

$$N = 2$$

$$k = 5$$

ordered frequent itemset

$$k = 5$$

$$y = 3$$

$$E = 4$$

$$m = 3$$

$$o = 3$$

} min-sup



E - 4  
 Y - 3  
 D - 1  
 A - 1  
 U - 1  
 C - 2  
 I - 1

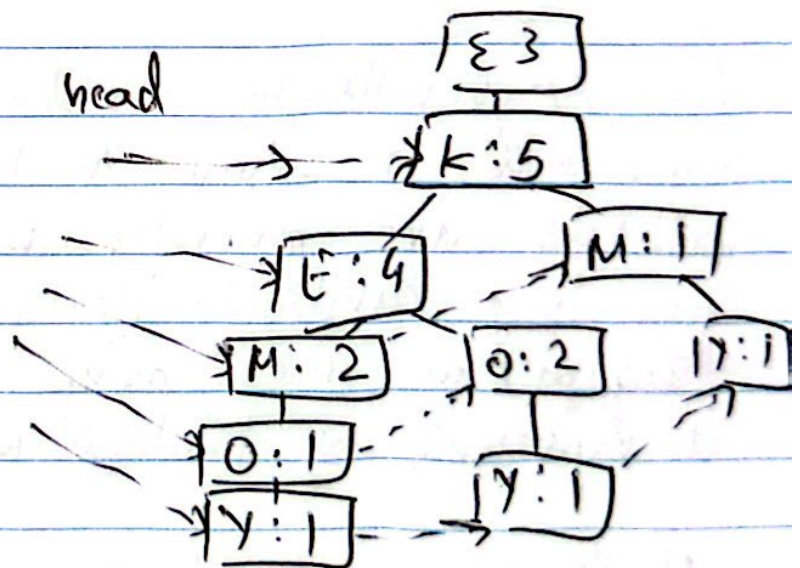
f list : k - E - M - O - Y

TID	Items - bought	ordered frequent item
T100	{M, O, N, K, E, Y}	{K, E, M, O, Y}
T200	{D, O, N, K, E, Y}	{K, E, O, Y}
T300	{M, A, K, E}	{K, E, M}
T400	{M, U, C, K, Y}	{K, M, Y}
T500	{C, O, O, K, I, E}	{K, E, O}

FP growth tree

Header table

Item	frequency	head
K	5	
E	4	
M	3	
O	3	
Y	3	





Items	Conditional Pattern Base	Conditional FPT	Frequent Pattern generated
Y	$\{ \{k, E, M, O: 1\} \}$ $\{k, E, O: 1\}, \{k, M: 1\}$	$\{k: 3\}$	$\{k, Y: 3\}$
O	$\{ \{k, E, M: 1\}, \{k, E: 2\} \}$	$\{k: 3\}, \{E: 3\}$	$\{k, E, O: 3\}, \{k, O: 3\}$ $\{E, O: 3\}$
M	$\{ \{k, E: 2\}, \{k: 1\} \}$	$\{k: 3\}$	$\{k, M: 3\}$
E	$\{ \{k: 4\} \}$	$\{k: 4\}$	$\{k, E: 4\}$

Frequent patterns generated using the F-P growth tree

fp growth tree  
 $\{ \{k, Y: 3\}, \{O, k: 3\}, \{O, E: 3\}, \{O, k, E: 3\}, \{M, k: 3\}, \{E, k: 4\}, \{E: 4\}, \{k: 4\}, \{M: 3\}, \{Y: 3\}, \{O: 3\} \}$

Apriori algorithm generates candidate itemsets may be large is no. if the dataset in the database is huge, it needs multiple scans of the database and this makes very slow...

In FP growth the problems in the Apriori algorithm is solved. In this algorithm frequent patterns are generated without the need of generating candidate itemsets. FP growth algorithm need only two db scans and it represents the database in the form of fp tree.

That's the reason fp growth algorithm is more efficient than a Apriori Algorithm



b) Meta Rule :  $\forall x \in \text{transaction},$   
 $\text{buys}(x, \text{item}_1) \wedge \text{buys}(x, \text{item}_2) \Rightarrow$   
 $\text{buys}(x, \text{item}_3)$

$x \rightarrow$  variable representing customers,  
 $\text{item}_i \rightarrow$  denotes variables representing items  
 (eg:- "A", "B", --- etc)

Association rules based on the above meta rule

$\forall x \in \text{transaction}, \text{buys}(x, E) \wedge \text{buys}(x, O) \Rightarrow \text{buys}(x, E)$   
 $\forall x \in \text{transaction}, \text{buys}(x, E) \wedge \text{buys}(x, O) \Rightarrow \text{buys}(x, E)$

3. (c) At the granularity of item-category (eg: milk give the rules of the template)

The largest frequent itemset for largest values of  $k=3$  is {milk, Bread, cheese}

①  $\text{Buys}(x, \text{milk}) \wedge \text{buys}(x, \text{Bread}) \Rightarrow \text{buys}(x, \text{cheese})$

$$\text{Support} = \frac{\sigma(\text{Milk, Bread, cheese})}{\text{Total transactions}} = \frac{3}{4} \times 100 = 75\%$$

$$\text{Confidence} = \frac{\sigma(\text{Milk, Bread, cheese})}{\sigma(\text{cheese})} = \frac{3}{3} \times 100 = 100\%$$

[75%, 100%]



$$\textcircled{2} \text{buys}(X, \text{milk}) \wedge \text{buys}(X, \text{cheese}) \Rightarrow \text{buys}(X, \text{Bread})$$

$$\text{Support} = \frac{\sigma(\text{milk}, \text{cheese}, \text{bread})}{\text{total transactions}(T)} = \frac{3}{4} \times 100 = 75\%.$$

[75%, 100%]

$$\textcircled{3} \text{Buys}(X, \text{Bread}) \wedge \text{buys}(X, \text{cheese}) \Rightarrow \text{buys}(X, \text{Milk})$$

$$\text{Support} = \frac{\sigma(\text{Bread}, \text{cheese}, \text{Milk})}{\text{total transactions}(T)} = \frac{3}{4} \times 100 = 75\%.$$

$$\text{Confidence} = \frac{\sigma(\text{Bread}, \text{cheese}, \text{Milk})}{\sigma(\text{Bread}, \text{cheese})} = \frac{3}{3} \times 100 = 100\%.$$

[75%, 100%]

b) The largest frequent itemset for  $k=3$  is  
 $\{ \text{wonder\_bread}, \text{sunset\_milk}, \text{Dairyland\_cheese} \}$   
 $\{ \text{Dairyland\_Milk}, \text{Tasty-pie}, \text{wonder\_Bread} \}$

4. For the mine association rules algorithm is as follows:-

→ firstly we need to find the local frequent itemsets in each store and here suppose,  $|F|$  be the union of all local frequent itemsets in the four stores, so in each store we need to find the local absolute support for each itemset in  $|F|$



And then we must determine the global (absolute) support for each itemset in IFI, this can be done by summing up for each itemset, the local support of that itemset in the four stores. Doing this for each itemset in IFI will give us their global supports. Itemset whose global supports pass the support threshold are good frequent itemsets.

Then lastly need to also derive strong association rules from the global frequent itemsets.

Thus there are the rules for the mine global association rules in an algorithm.

50 Given

	hotdogs	hotdogs(NOT)	total
Hamburgers	2000	500	2500
Hamburgers NOT	1000	1500	2500
total	3000	2000	5000

(a) Given rule :

hotdogs  $\Rightarrow$  hamburgers  
 $\text{min\_sup} = 25\%$



$$\text{min-confidence} = 50\%$$

$$\text{Support} = \frac{2000}{5000} = 40\%$$

$$\text{confidence} = \frac{2000}{5000} = 66.7\%$$

Support 40%, confidence 66.7% is greater than the given minimum support and minimum confidence.

∴ The given rule association is strong

$$\begin{aligned} \text{b) Correlation (hotdog, hamburger)} \\ &= \frac{P(\text{hotdog, hamburger})}{P(\text{hotdog}) \cdot P(\text{hamburger})} \end{aligned}$$

$$\text{Cor}(A, B) = \frac{P(A \cap B)}{P(A) \cdot P(B)} = \frac{(2000/5000)}{(\frac{3000}{5000})(\frac{2500}{5000})}$$

$$> \frac{0.4}{0.6 \times 0.5} > 1.33 (> 1)$$

So the correlation value is greater than 1  
Therefore the purchase of hotdog is not independent of the purchase of hamburger and there is a positive correlation.