

09-09-2020

Naman Garg

B032

BTech CS B.

## DATA MINING ASSIGNMENT 2

Apriori Algorithm :  $\text{Support} = 60\% = \frac{60}{100} \times 5 = \underline{3}$

Q1 (i) Count number of support for each item  $C_i$

$C_i$	Support Count
A	1
C	2
D	1
E	4
I	1
K	5
M	3
N	2
O	3
U	1
Y	3

FREQUENT ITEMSET  $\text{Support} \geq 3$

$L_i$	Support Count
E	4
K	5
M	3
O	3
Y	3

(ii) SELF JOIN:  $L_1 \rightarrow C_2$   $L_1 \bowtie L_1$

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

$C_2$	Support count
E, K	4
E, M	2
E, O	3
E, Y	2
K, M	3
K, O	3
K, Y	3
M, O	1
M, Y	2
O, Y	2

Frequent item set  $L_2$   
Support Count  $\geq 3$

$L_2$	Support Count
E, K	4
E, O	3
K, M	3
K, O	3
K, Y	3

(iii) Self join  $L_2 \rightarrow C_3$  using the prune technique  
i.e. we must check the fact that all non-empty subsets of a frequent itemset must be a frequent itemset as well.

$$L_2 = \{EK, EO, KM, KO, KY\} \cap \{EK, EO, KM, KO, KY\} \\ = \{EKO, EKM, EKY, KMO, KMY\}$$

$EKM \Rightarrow \{E, K\}$  and  $\{K, M\}$  are frequent item sets but  $\{E, M\}$  is not ~~here~~ hence, we remove it.

$EKY \Rightarrow \{E, Y\}$  is not frequent item set, we remove it.

all subsets of:

$\{EKO\}$ ,  $\{KMY\}$  and  $\{KMO\}$  are included as frequent subsets  $\therefore$  we keep them.

thus, we have:

$C_3$	Support Count
$EKO$	3
$KMO$	1
$KMY$	2

 $\therefore L_3 \Rightarrow$ 

$L_3$	Count
$EKO$	3

There is only 1 value in  $L_3$   $\therefore$  we cannot make an  $C_4$ .

$\therefore$  No more associations are possible.

### ASSOCIATION RULES

- $[K, O] \rightarrow E = 3/3 = 100\%$
- $[E, K] \rightarrow O = 3/4 = 75\%$
- $[E, O] \rightarrow K = 3/3 = 100\%$
- $E \rightarrow [K, O] = 3/4 = 75\%$
- $O \rightarrow [E, K] = 3/3 = 100\%$
- $K \rightarrow [E, O] = 3/5 = 60\%$

we want confidence  $\geq 80\%$ .

$\therefore$  we discard ~~2, 4, 6~~ we keep 1, 3, 5  
because they have confidence  $\geq 80\%$



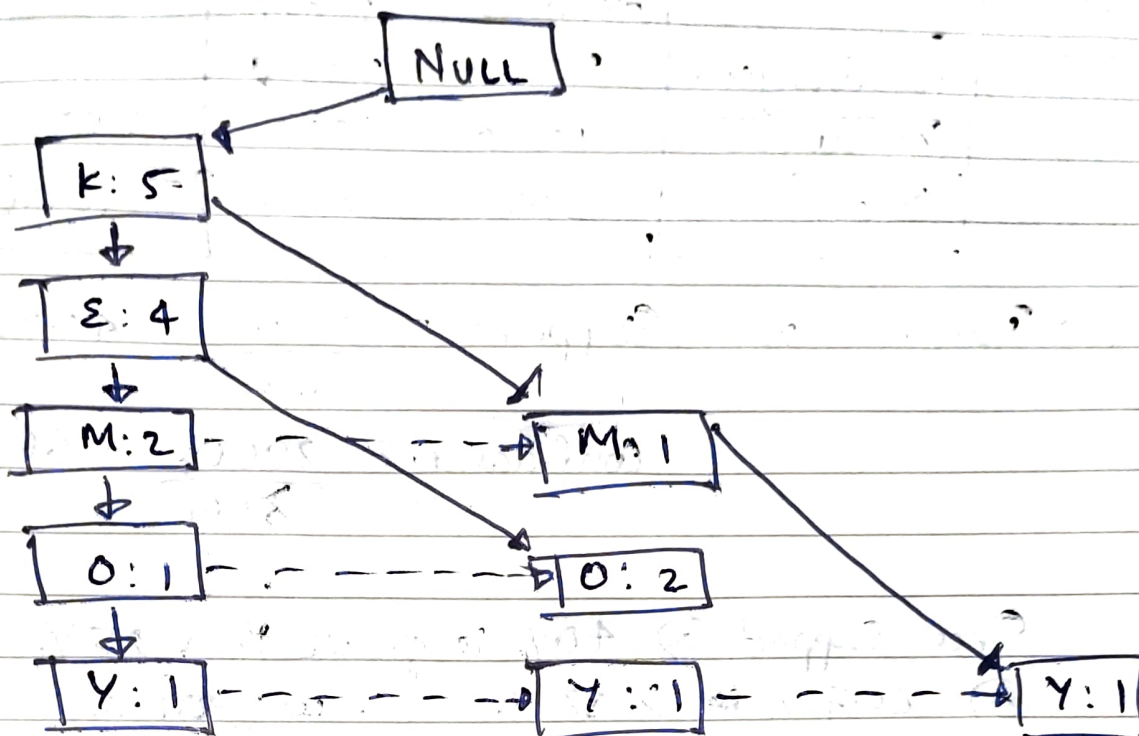
Using FP Growth: Min Support = 60%  $\frac{60}{100} \times 5 = 3$

Item	Frequency
A	1
C	2
D	1
E	4
I	1
K	5
M	3
N	2
O	3
U	1
Y	3

2. Frequent Pattern set  $L = \{K:5, E:4, M:3, O:3, Y:3\}$

Transaction ID	Items	Ordered-Item set
T100	{M, O, N, E, Y}	{E, O, M, Y}
T200	{D, O, N, E, Y}	{E, O, Y}
T300	{M, A, K, E}	{E, M}
T400	{M, U, C, Y}	{K, M, Y}
T500	{C, O, U, K, I, E}	{E, O, Y}

## FP Tree



Items	Frequent Pattern Generated
Y	$\{ \langle K, Y: 3 \rangle \}$
O	$\{ \langle K, O: 3 \rangle, \langle E, O: 3 \rangle, \langle E, K, O: 3 \rangle \}$
M	$\{ \langle K, M: 3 \rangle \}$
E	$\{ \langle E, K: 3 \rangle \}$
K	

## Association Rules

- $[E, K] \rightarrow O = 3/4 = 75\%$
- $[K, O] \rightarrow E = 3/3 = 100\%$
- $[E, O] \rightarrow K = 3/3 = 100\%$
- $E \rightarrow [K, O] = 3/4 = 75\%$
- $K \rightarrow [E, O] = 3/5 = 60\%$
- $O \rightarrow [E, K] = 3/3 = 100\%$

Rules 1, 4, 5 are discarded because we want confidence  $\geq 80\%$ . 2, 3, 6 are selected.

Q2

	Hotdogs	Hottogs	$\Sigma$ row
hamburgers	2000	3000	2500
<u>hamburgers</u>	1000	1500	2500
$\Sigma$ col.	3000	2000	5000

(i) For the rule,  $\text{Support} = \frac{2000}{5000} = 40\%$ .

$$\text{Confidence} = \frac{2000}{3000} = 66.67\%$$

Since Support of 40% is greater than 25% & confidence of 66.7% > 50%.  $\therefore$  the association rule is STRONG

$$\begin{aligned} \text{(ii) Corr (hotdog, hamburger)} &= \frac{P(\{\text{hot dog, hamburger}\})}{(P(\{\text{hot dog}\}) P(\{\text{hamburger}\}))} \\ &= \frac{0.4}{(0.5 \times 0.6)} = 1.33 > 1 \end{aligned}$$

$\therefore$  Since the Corr value > 1 the purchase of hotdogs is NOT independent of the purchase of hamburgers. There ~~are~~ is a ~~are~~ POSITIVE correlation b/w them.