

NAME: MD EZAZ ANWAR

ROLL NO.: IIT2018004

COURSE: DATA MINING

SEMESTER: VI

ASSIGNMENT - 2

1. Apply apriori algorithm on the following datasets.

A.

TID	Items
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I1, I2, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I2, I3

(i) Support = 20%, Confidence = 20%.

Given, support count =  $\frac{20}{100} \times 9 = 1.8 = 2$ 

	C <sub>1</sub>		L <sub>1</sub>		
	Itemset	Sup. Count	Compare candidate	Itemset	Sup. Count
Scan D for count of each candidate	{I1}	6	support count with minimum	{I1}	6
	{I2}	7		{I2}	7
	{I3}	6	support count	{I3}	6
	{I4}	2		{I4}	2
→	{I5}	2		{I5}	2

$C_2$			$L_2$		
	Itemset	Sup. Count		Itemset	Sup. Count
Generate $C_2$ candidates from $L_1$ and scan D for count of each candidate	$\{I_1, I_2\}$	4	Compare candidate support count with minimum support count	$\{I_1, I_2\}$	4
	$\{I_1, I_3\}$	4		$\{I_1, I_3\}$	4
	$\{I_1, I_4\}$	1		$\{I_1, I_5\}$	2
	$\{I_1, I_5\}$	2		$\{I_2, I_3\}$	4
	$\{I_2, I_3\}$	4		$\{I_2, I_4\}$	2
	$\{I_2, I_4\}$	2		$\{I_2, I_5\}$	2
→	$\{I_2, I_5\}$	2			
	$\{I_3, I_4\}$	0			
	$\{I_3, I_5\}$	1			
	$\{I_4, I_5\}$	0			

$C_3$			$L_3$		
	Itemset	Sup. Count		Itemset	Sup. Count
Generate $C_3$ candidates from $L_2$ and scan D for count of each candidate	$\{I_1, I_2, I_3\}$	2	Compare candidate sup. count with min. sup. count	$\{I_4, I_2, I_3\}$	2
→	$\{I_1, I_2, I_5\}$	2	→	$\{I_1, I_2, I_5\}$	2

When we generate  $C_4$  using  $L_3 \bowtie L_3$ , we get  $\{\{I_1, I_2, I_3, I_5\}\}$ , where itemset  $\{I_1, I_2, I_3, I_5\}$  is pruned because its subset  $\{I_2, I_3, I_5\}$  is not frequent. Thus,  $C_4 = \emptyset$  and algorithm terminates.

Now, we have frequent itemsets of size greater than 2,

$$L = \{\{I_1, I_2\}, \{I_1, I_3\}, \{I_1, I_5\}, \{I_2, I_3\}, \{I_2, I_4\}, \{I_2, I_5\}, \{I_1, I_2, I_3\}, \\ \{I_1, I_2, I_5\}\}$$

So, the association rules we get are,

$$\text{I. } \{I_1\} \Rightarrow \{I_2\} \quad \text{confidence} = 4/6 = 67\%$$

$$\text{II. } \{I_2\} \Rightarrow \{I_1\} \quad \text{confidence} = 4/7 = 57\%$$

$$\text{III. } \{I_1\} \Rightarrow \{I_3\} \quad \text{confidence} = 4/6 = 67\%$$

$$\text{IV. } \{I_3\} \Rightarrow \{I_1\} \quad \text{confidence} = 4/6 = 67\%$$

V.	$\{I_1\} \Rightarrow \{I_5\}$	confidence = $2/6 = 33\%$ .
VI.	$\{I_5\} \Rightarrow \{I_1\}$	confidence = $2/2 = 100\%$ .
VII.	$\{I_2\} \Rightarrow \{I_3\}$	confidence = $4/7 = 57\%$ .
VIII.	$\{I_3\} \Rightarrow \{I_2\}$	confidence = $4/6 = 67\%$ .
IX.	$\{I_2\} \Rightarrow \{I_4\}$	confidence = $2/7 = 29\%$ .
X.	$\{I_4\} \Rightarrow \{I_2\}$	confidence = $2/2 = 100\%$ .
XI.	$\{I_2\} \Rightarrow \{I_5\}$	confidence = $2/7 = 29\%$ .
XII.	$\{I_5\} \Rightarrow \{I_2\}$	confidence = $2/2 = 100\%$ .
XIII.	$\{I_1, I_2\} \Rightarrow \{I_3\}$	confidence = $2/4 = 50\%$ .
XIV.	$\{I_1, I_3\} \Rightarrow \{I_2\}$	confidence = $2/4 = 50\%$ .
XV.	$\{I_2, I_3\} \Rightarrow \{I_1\}$	confidence = $2/4 = 50\%$ .
XVI.	$\{I_3\} \Rightarrow \{I_1, I_2\}$	confidence = $2/6 = 33\%$ .
XVII.	$\{I_2\} \Rightarrow \{I_1, I_3\}$	confidence = $2/7 = 29\%$ .
XVIII.	$\{I_1\} \Rightarrow \{I_2, I_3\}$	confidence = $2/6 = 33\%$ .
XIX.	$\{I_1, I_2\} \Rightarrow \{I_5\}$	confidence = $2/4 = 50\%$ .
XX.	$\{I_1, I_5\} \Rightarrow \{I_2\}$	confidence = $2/2 = 100\%$ .
XXI.	$\{I_2, I_5\} \Rightarrow \{I_1\}$	confidence = $2/2 = 100\%$ .
XXII.	$\{I_5\} \Rightarrow \{I_1, I_2\}$	confidence = $2/2 = 100\%$ .
XXIII.	$\{I_2\} \Rightarrow \{I_1, I_5\}$	confidence = $2/7 = 29\%$ .
XXIV.	$\{I_1\} \Rightarrow \{I_2, I_5\}$	confidence = $2/6 = 33\%$ .

All the association rules have confidence  $\geq$  minimum confidence. Hence, all the rules are strong association rules.

(ii) Support = 30%, Confidence = 30%.

Given, support count =  $\frac{30}{100} \times 9 = 2.7 \approx 3$

(a) Scan D for count of each candidate,

C <sub>1</sub>	Itemset	{I <sub>1</sub> }	{I <sub>2</sub> }	{I <sub>3</sub> }	{I <sub>4</sub> }	{I <sub>5</sub> }
	Sup. Count	6	7	6	2	2

Compare Support

count with

minimum support

count

Itemset	Sup. Count
$\{I_1\}$	6
$\{I_2\}$	7
$\{I_3\}$	6

Generate  $C_2$

from  $L_1$  and

scan D for count

$C_2$  combinations

Itemset	Count
$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_2, I_3\}$	4

Compare cand.

sup. count with

min. sup. count

$L_2$

Itemset	Count
$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_2, I_3\}$	4

Generate  $C_3$  from  $L_2$   
and scan D for count

Itemset	Sup. Count
$\{I_1, I_2, I_3\}$	2

when we generate  $L_3$  from  $C_3$  we get  $L_3 = \emptyset$  and algorithm terminates.

So, we have,  $L = \{\{I_1, I_2\}, \{I_1, I_3\}, \{I_2, I_3\}\}$

Hence, association rules,

I.  $\{I_1\} \Rightarrow \{I_2\}$  confidence =  $4/6 = 67\%$

II.  $\{I_2\} \Rightarrow \{I_1\}$  confidence =  $4/7 = 57\%$

III.  $\{I_1\} \Rightarrow \{I_3\}$  confidence =  $4/6 = 67\%$

IV.  $\{I_3\} \Rightarrow \{I_1\}$  confidence =  $4/6 = 67\%$

V.  $\{I_2\} \Rightarrow \{I_3\}$  confidence =  $4/7 = 57\%$

VI.  $\{I_3\} \Rightarrow \{I_2\}$  confidence =  $4/6 = 67\%$

All the association rules have confidence  $\geq 30\%$ . Hence, all rules are strong association rules.

B.	Transaction ID	Items
	T1	{E, K, M, N, O, Y}
	T2	{D, E, K, N, O, Y}
	T3	{A, E, K, M}
	T4	{C, K, M, U, Y}
	T5	{C, E, I, K, O, O}

(i) Support = 40%, Confidence = 40%.

$$\text{we have, support count} = \frac{40}{100} \times 5 = 2$$

STEP I: K=1

Scan D for count of each candidate.

C <sub>1</sub>	Itemset	{E}	{K}	{M}	{N}	{O}	{Y}	{D}	{A}	{C}	{U}	{I}
	Count	4	5	3	2	3	3	1	1	2	1	1

Compare candidate support count with minimum support count.

L <sub>1</sub>	Itemset	{E}	{K}	{M}	{N}	{O}	{Y}	{C}
	Count	4	5	3	2	3	3	2

STEP II: K=2

Generate C<sub>2</sub> using L<sub>1</sub> and count each candidate by scanning D.

C <sub>2</sub>	Itemset	{E, K}	{E, M}	{E, N}	{E, O}	{E, Y}	{E, C}	{K, M}	{K, N}	{K, O}	{K, Y}	{K, C}
	Count	4	2	2	3	2	1	2	2	3	3	2

C <sub>2</sub>	Itemset	{M, N}	{M, O}	{M, Y}	{M, C}	{N, O}	{N, Y}	{N, C}	{O, Y}	{O, C}	{Y, C}
	Count	1	1	2	1	2	2	0	2	1	1

Compare candidate support count with minimum support count.

$L_2$	Itemset	{E, K}	{E, M}	{E, N}	{E, O}	{E, Y}	{K, M}	{K, N}
	Count	4	2	2	3	2	2	2

$L_2$	Itemset	{K, O}	{K, Y}	{K, C}	{M, Y}	{N, O}	{N, Y}	{O, Y}
	Count	3	3	2	2	2	2	2

STEP - III :  $K = 3$

Generate  $C_3$  using  $L_2$  and scan D for count of each candidate.

$C_3$	Itemset	{E, K, M}	{E, K, N}	{E, K, O}	{E, K, Y}	{E, M, Y}	{E, N, O}	{E, N, Y}
	Count	2	2	3	2	1	2	2

$C_3$	Itemset	{E, O, Y}	{K, M, Y}	{K, N, O}	{K, N, Y}	{K, O, Y}	{N, O, Y}	
	Count	2	2	2	2	2	2	

Compare candidate support count with minimum support count.

$L_3$	Itemset	{E, K, M}	{E, K, N}	{E, K, O}	{E, K, Y}	{E, N, O}	{E, N, Y}	
	Count	2	2	3	2	2	2	

$L_3$	Itemset	{E, O, Y}	{K, M, Y}	{K, N, O}	{K, N, Y}	{K, O, Y}	{N, O, Y}	
	Count	2	2	2	2	2	2	

STEP - IV :  $K = 4$

Generate  $C_4$  using  $L_3$  and scan D for count of each candidate.

$C_4$	Itemset	{E, N, O, Y}	{E, K, N, O}	{E, K, O, Y}	{K, N, O, Y}	{E, K, N, Y}		
	Count	2	2	2	2	2		

Compare candidate support count with minimum support count.

$L_4$	Itemset	{E, N, O, Y}	{E, K, N, O}	{E, K, O, Y}	{K, N, O, Y}	{E, K, N, Y}		
	Count	2	2	2	2	2		

STEP-V:  $K=5$

Generate  $C_5$  using  $L_4$  and scan D for count of each candidate.

$C_5$

Itemset	Count
{E, K, N, O, Y}	2

Compare candidate support count with minimum sup. count.

$L_5$

Itemset	Count
{E, K, N, O, Y}	2

When we generate  $C_6 = L_5 \bowtie L_5$  we get  $C_6 = \emptyset$ , hence algorithm terminates.

So, we have frequent itemset,

$$L = L_2 \cup L_3 \cup L_4 \cup L_5$$

There are 14 elements in  $L_2$ , so it will give  $14 \times {}^2C_1 = 28$  rules

Similarly  $L_3$  will give  $12 \times {}^3C_1 \times 2 = 72$  rules

And,  $L_4$  will give  $5 \times ({}^4C_1 \times 2 + {}^4C_2) = 70$  rules

And,  $L_5$  will give  $1 \times ({}^5C_1 \times 2 + {}^5C_2) = 30$  rules

i.e. a total of 200 association rules. Hence, only some rules are shown below.

using, {E, K} we get two rules

I.  $\{E\} \Rightarrow \{K\}$  confidence =  $4/4 = 100\%$ .

II.  $\{K\} \Rightarrow \{E\}$  confidence =  $4/5 = 80\%$ .

using, {E, K, O} we get six rules

III.  $\{E, K\} \Rightarrow \{O\}$  confidence =  $3/4 = 75\%$ .

IV.  $\{E, O\} \Rightarrow \{K\}$  confidence =  $3/3 = 100\%$ .

V.  $\{K, O\} \Rightarrow \{E\}$  confidence =  $3/3 = 100\%$ .

VI.  $\{O\} \Rightarrow \{E, K\}$  confidence =  $3/3 = 100\%$ .

VII.  $\{K\} \Rightarrow \{E, O\}$  confidence =  $3/5 = 60\%$ .

VIII.  $\{E\} \Rightarrow \{K, O\}$  confidence =  $3/4 = 75\%$ .

Using  $\{E, K, O, Y\}$ , we will have 14 rules, some of which are,

IX.  $\{E, K, O\} \Rightarrow \{Y\}$  confidence =  $2/3 = 67\%$ .

X.  $\{K\} \Rightarrow \{E, O, Y\}$  confidence =  $2/5 = 40\%$ .

XI.  $\{E, O\} \Rightarrow \{K, Y\}$  confidence =  $2/3 = 67\%$ .

Using  $\{E, K, N, O, Y\}$ , we will have 30 rules, some of which are,

XII.  $\{E, K, N, O\} \Rightarrow \{Y\}$  confidence =  $2/2 = 100\%$ .

XIII.  $\{E, K, N\} \Rightarrow \{O, Y\}$  confidence =  $2/2 = 100\%$ .

XIV.  $\{E, K\} \Rightarrow \{N, O, Y\}$  confidence =  $2/4 = 50\%$ .

XV.  $\{E\} \Rightarrow \{K, N, O, Y\}$  confidence =  $2/4 = 50\%$ .

All the rules have confidence  $\geq$  min-confidence. So, all rules are strong association rules.

(ii). Support = 3 Confidence = 20%.

Given, support = 3

STEP-I:  $K=1$

Scan D to find the count of each candidate.

C <sub>1</sub>	Itemset	{E}	{K}	{M}	{N}	{O}	{Y}	{D}	{A}	{C}	{U}	{J}
Count	4	5	3	2	3	3	1	1	1	2	1	1

Compare candidate count with minimum support count.

L <sub>1</sub>	Itemset	{E}	{K}	{M}	{O}	{Y}
Count	4	5	3	3	3	AT

STEP-II:  $K=2$

Scan D to find the count of each candidate generated in C<sub>2</sub>.

C <sub>2</sub>	Itemset	{E, K}	{E, M}	{E, O}	{E, Y}	{K, M}	{K, O}	{K, Y}	{M, O}	{M, Y}	{O, Y}
Count	4	2	3	2	3	3	3	1	2	2	

Compare candidate count with minimum support count.

L <sub>2</sub>	Itemset	{E, K}	{E, O}	{K, M}	{K, O}	{K, Y}
Count	4	3	3	3	3	

STEP - III :  $K=3$

generate  $C_3$  using  $L_2$  and scan  $D$  to find the count of each cand.

$C_3$	
Itemset	Count
{E, K, O}	3

compare candidate support count with min. support count.

$C_3$	
Itemset	Count
{E, K, O}	3

when we generate  $C_4$  using  $L_3$ , we get  $C_4 = \emptyset$  and algorithm terminates.

Now, we have,

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

so, the association rules are,

- i.  $\{E\} \Rightarrow \{K\}$  confidence =  $4/4 = 100\%$ .
- ii.  $\{K\} \Rightarrow \{E\}$  confidence =  $4/5 = 80\%$ .
- iii.  $\{E\} \Rightarrow \{O\}$  confidence =  $3/4 = 75\%$ .
- iv.  $\{O\} \Rightarrow \{E\}$  confidence =  $3/3 = 100\%$ .
- v.  $\{K\} \Rightarrow \{M\}$  confidence =  $3/5 = 60\%$ .
- vi.  $\{M\} \Rightarrow \{K\}$  confidence =  $3/3 = 100\%$ .
- vii.  $\{K\} \Rightarrow \{O\}$  confidence =  $3/5 = 60\%$ .
- viii.  $\{O\} \Rightarrow \{K\}$  confidence =  $3/3 = 100\%$ .
- ix.  $\{K\} \Rightarrow \{Y\}$  confidence =  $3/5 = 60\%$ .
- x.  $\{Y\} \Rightarrow \{K\}$  confidence =  $3/3 = 100\%$ .
- xi.  $\{E, K\} \Rightarrow \{O\}$  confidence =  $3/4 = 75\%$ .
- xii.  $\{O\} \Rightarrow \{E, K\}$  confidence =  $3/3 = 100\%$ .
- xiii.  $\{E, O\} \Rightarrow \{K\}$  confidence =  $3/3 = 100\%$ .
- xiv.  $\{K\} \Rightarrow \{E, O\}$  confidence =  $3/5 = 60\%$ .
- xv.  $\{O, K\} \Rightarrow \{E\}$  confidence =  $3/3 = 100\%$ .
- xvi.  $\{E\} \Rightarrow \{O, K\}$  confidence =  $3/4 = 75\%$ .

All the above rules are strong association rules, since every rule has  $conf \geq 20\%$ .

2. Apply FP-growth algorithm on the following datasets.

A.

Transaction ID	Items
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I1, I2, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I2, I3

support = 20% and confidence = 20%

Given, support count =  $0.2 \times 9 = 1.8 \approx 2$

scan the database to find the count of each item.

Item	I1	I2	I3	I4	I5
count	6	7	6	2	2

The set L is constructed in the order of decreasing support count

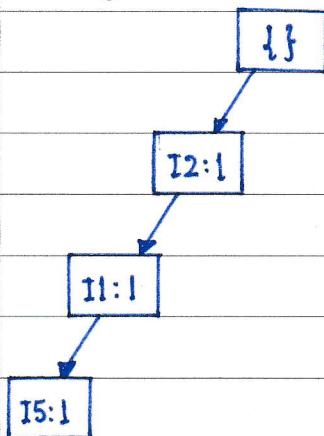
$$L = \{\{I2: 7\}, \{I1: 6\}, \{I3: 6\}, \{I4: 2\}, \{I5: 2\}\}$$

Now, ordered itemset is built by scanning database and L.

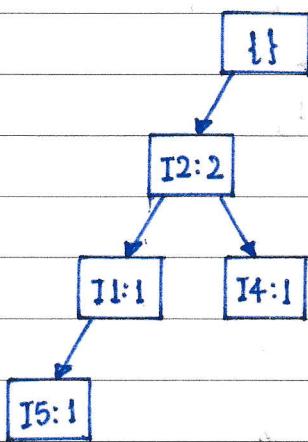
Transaction ID	Ordered Itemset	Transaction ID	Ordered Itemset
T1	{I2, I1, I5}	T6	{I2, I3}
T2	{I2, I4}	T7	{I1, I3}
T3	{I2, I3}	T8	{I1, I2, I3, I5}
T4	{I2, I1, I4}	T9	{I2, I1, I3}
T5	{I1, I3}		

Now, we will make the FP-tree.

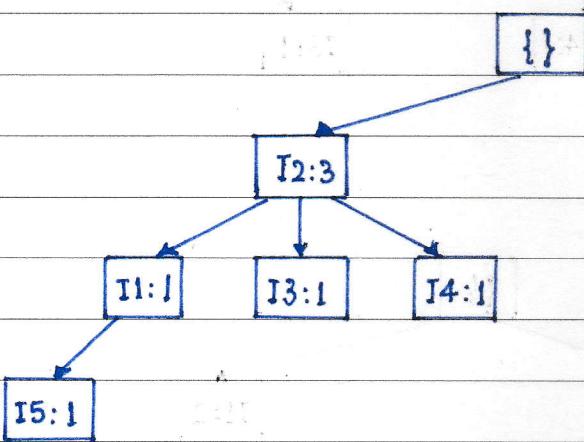
(i) Inserting {I2, I1, I5},



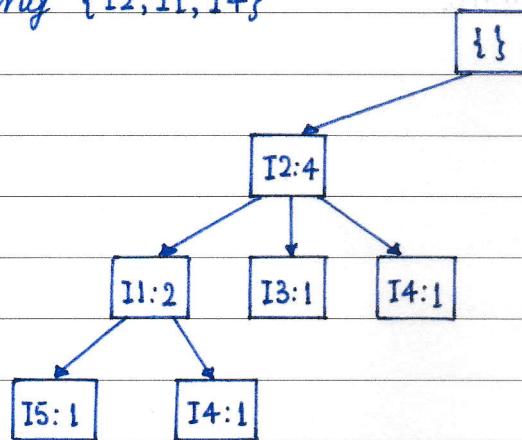
(ii) Inserting {I2, I4}



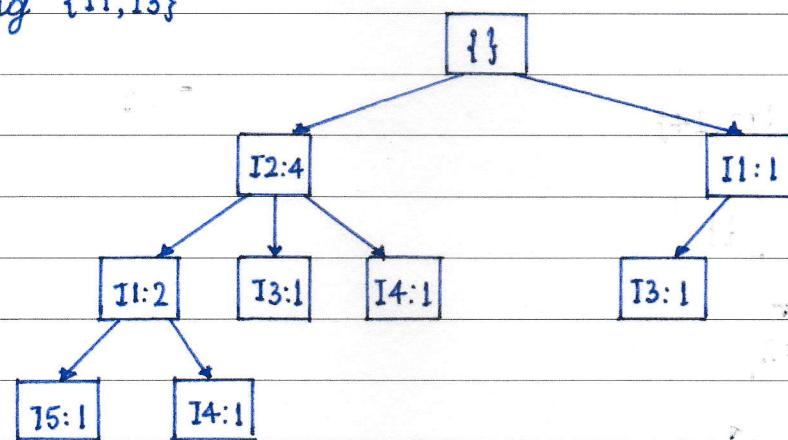
(iii) Inserting {I2, I3}



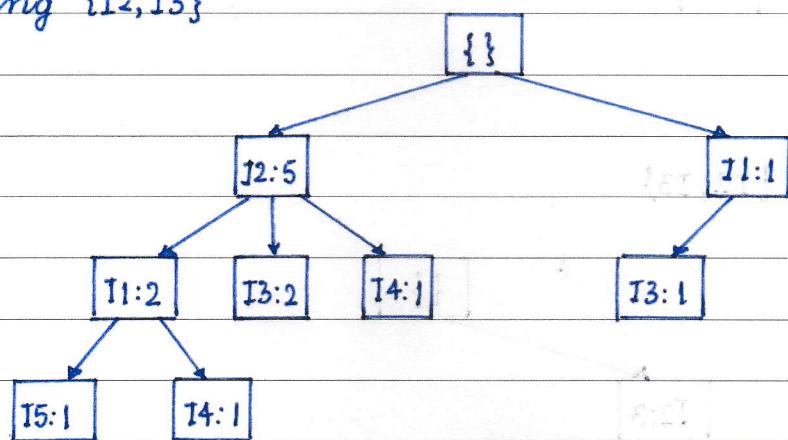
(iv) Inserting  $\{I_2, I_1, I_4\}$



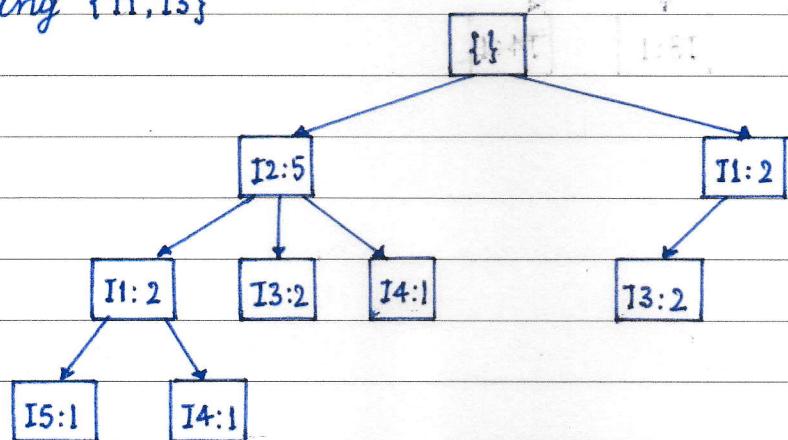
(v) Inserting  $\{I_1, I_3\}$



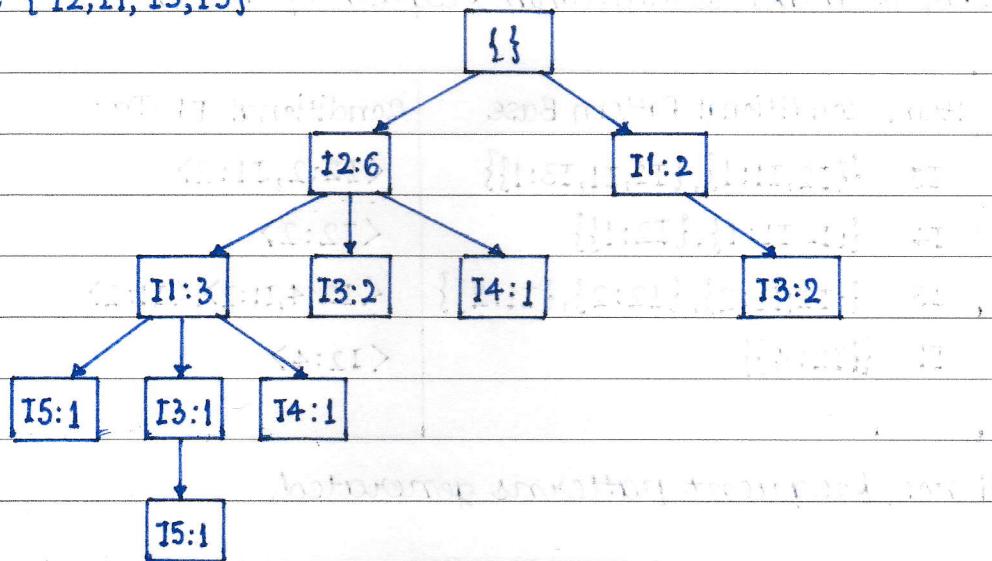
(vi) Inserting  $\{I_2, I_3\}$



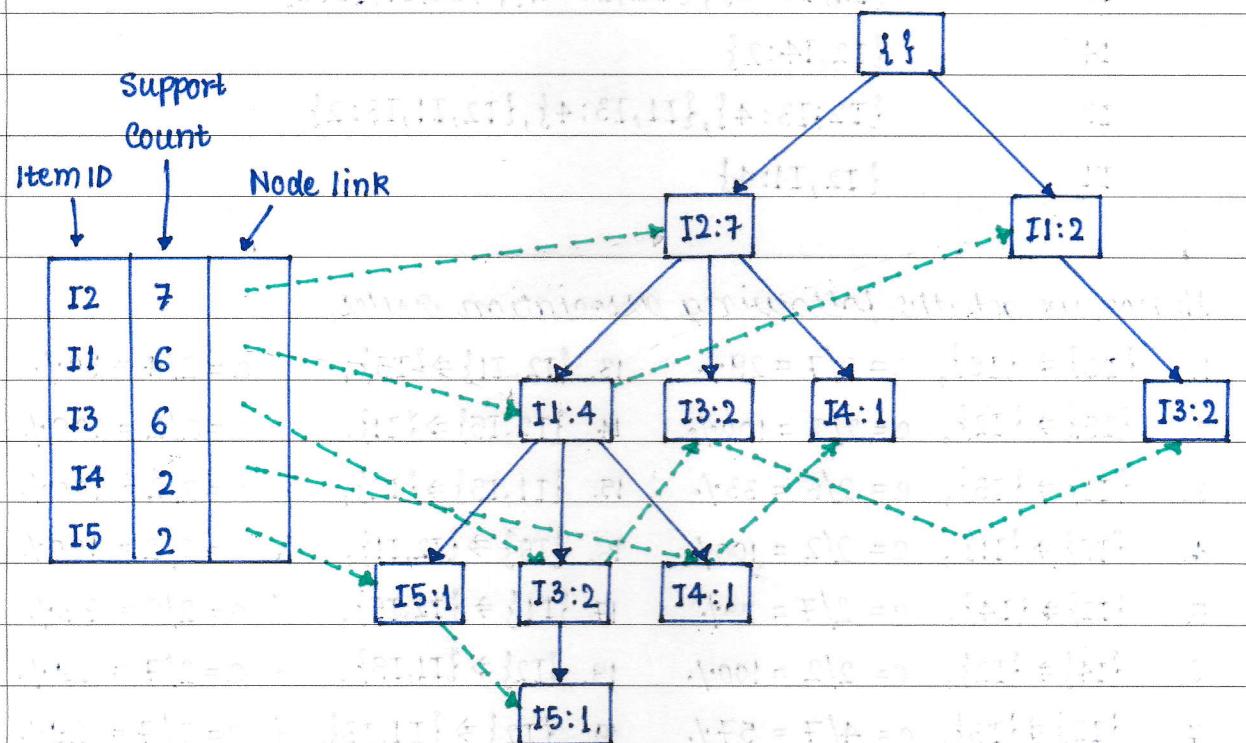
(vii) Inserting  $\{I_1, I_3\}$



(viii) insert {I2, I1, I3, I5}



(ix) After inserting {I2, I1, I3}, the final FP-tree,



Now, we generate Conditional Pattern Base:

Items	Conditional Pattern Base
I5	{I2, I1:1}, {I2, I1, I3:1}
I4	{I2, I1:1}, {I2:1}
I3	{I2, I1:2}, {I2:2}, {I1:2}
I1	{I2:4}

Now, we make conditional frequent pattern tree,

Items	Conditional Pattern Base	Conditional FP-Tree
I5	{I2, I1: 1}, {I2, I1, I3: 1}	<I2: 2, I1: 2>
I4	{I2, I1: 1}, {I2: 1}	<I2: 2>
I3	{I2; I1: 2}, {I2: 2}, {I1: 2}	<I2: 4, I1: 2>, <I1: 2>
I1	{I2: 4}	<I2: 4>

Hence, frequent patterns generated,

Items	Frequent Pattern Generated
I5	{I2, I5: 2}, {I2, I5: 2}, {I2, I1, I5: 2}
I4	{I2, I4: 2}
I3	{I2; I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2}
I1	{I2, I1: 4}

Hence, we get the following association rules,

- |  |  |
|--|--|
| 1. $\{I_2\} \Rightarrow \{I_5\}$ , $c = 2/7 = 29\%$ .  | 13. $\{I_2, I_1\} \Rightarrow \{I_5\}$ , $c = 2/4 = 50\%$ .  |
| 2. $\{I_5\} \Rightarrow \{I_2\}$ , $c = 2/2 = 100\%$ . | 14. $\{I_2, I_5\} \Rightarrow \{I_1\}$ , $c = 2/2 = 100\%$ . |
| 3. $\{I_1\} \Rightarrow \{I_5\}$ , $c = 2/6 = 33\%$ .  | 15. $\{I_1, I_5\} \Rightarrow \{I_2\}$ , $c = 2/2 = 100\%$ . |
| 4. $\{I_5\} \Rightarrow \{I_1\}$ , $c = 2/2 = 100\%$ . | 16. $\{I_5\} \Rightarrow \{I_2, I_1\}$ , $c = 2/2 = 100\%$ . |
| 5. $\{I_2\} \Rightarrow \{I_4\}$ , $c = 2/7 = 29\%$ .  | 17. $\{I_1\} \Rightarrow \{I_2, I_5\}$ , $c = 2/6 = 33\%$ .  |
| 6. $\{I_4\} \Rightarrow \{I_2\}$ , $c = 2/2 = 100\%$ . | 18. $\{I_2\} \Rightarrow \{I_1, I_5\}$ , $c = 2/7 = 29\%$ .  |
| 7. $\{I_2\} \Rightarrow \{I_3\}$ , $c = 4/7 = 57\%$ .  | 19. $\{I_2\} \Rightarrow \{I_1, I_3\}$ , $c = 2/7 = 29\%$ .  |
| 8. $\{I_3\} \Rightarrow \{I_2\}$ , $c = 4/6 = 67\%$ .  | 20. $\{I_1\} \Rightarrow \{I_2, I_3\}$ , $c = 2/6 = 33\%$ .  |
| 9. $\{I_1\} \Rightarrow \{I_3\}$ , $c = 4/6 = 67\%$ .  | 21. $\{I_3\} \Rightarrow \{I_1, I_2\}$ , $c = 2/6 = 33\%$ .  |
| 10. $\{I_3\} \Rightarrow \{I_1\}$ , $c = 4/6 = 67\%$ . | 22. $\{I_1, I_2\} \Rightarrow \{I_3\}$ , $c = 2/4 = 50\%$ .  |
| 11. $\{I_2\} \Rightarrow \{I_1\}$ , $c = 4/7 = 57\%$ . | 23. $\{I_1, I_3\} \Rightarrow \{I_2\}$ , $c = 2/4 = 50\%$ .  |
| 12. $\{I_1\} \Rightarrow \{I_2\}$ , $c = 4/6 = 67\%$ . | 24. $\{I_2, I_3\} \Rightarrow \{I_1\}$ , $c = 2/4 = 50\%$ .  |

All the association rules have confidence  $\geq 20\%$ . So, all rules are strong association rules.

B.	Transaction ID	Items
	T1	E, K, M, N, O, Y
	T2	D, E, K, N, O, Y
	T3	A, E, K, M
	T4	C, K, M, U, Y
	T5	C, E, I, K, O, O

support = 3 and confidence = 20%

scan the database to find the count of each item.

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

The set L is constructed in the order of decreasing support count.

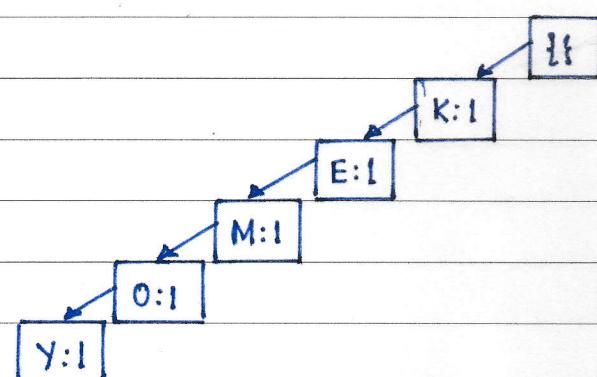
$$L = \{ \{K: 5\}, \{E: 4\}, \{M: 3\}, \{O: 3\}, \{Y: 3\} \}$$

Now, the ordered itemset is built by scanning the database,

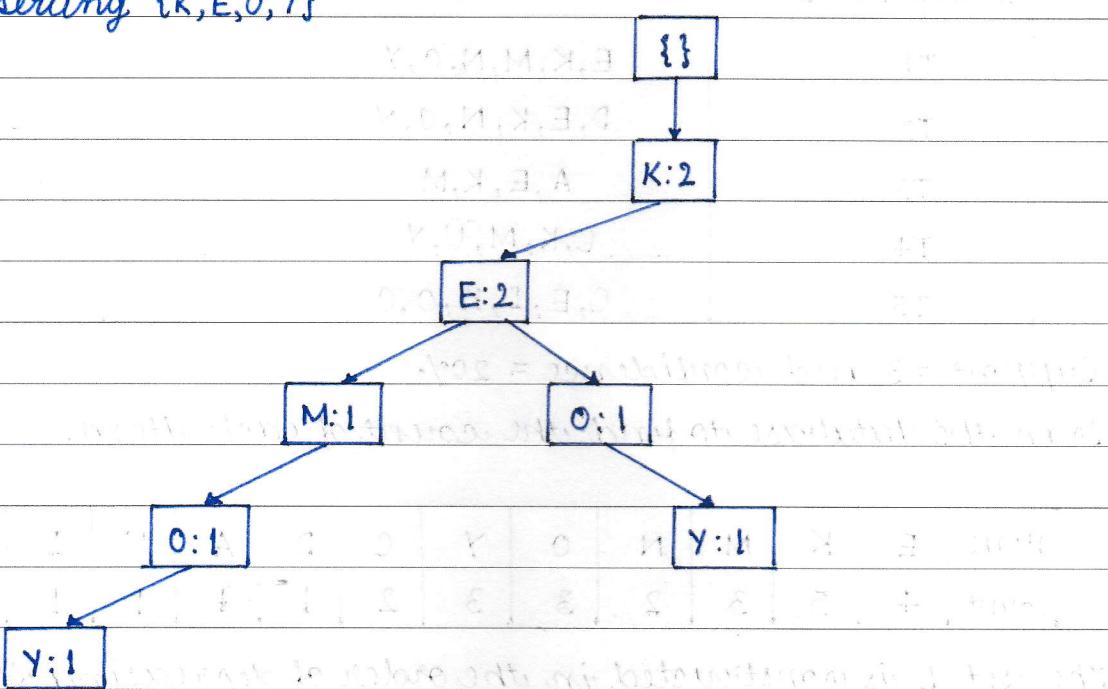
Transaction ID	Ordered Itemset
T1	{K, E, M, O, Y}
T2	{K, E, O, Y}
T3	{K, E, M}
T4	{K, M, Y}
T5	{K, E, O}

Now, we will make the FP-Tree

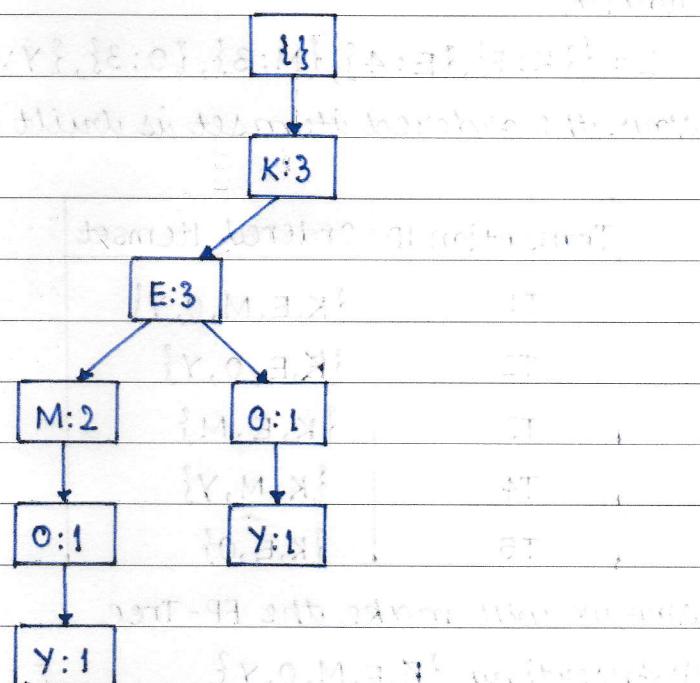
(i) Inserting {K, E, M, O, Y}



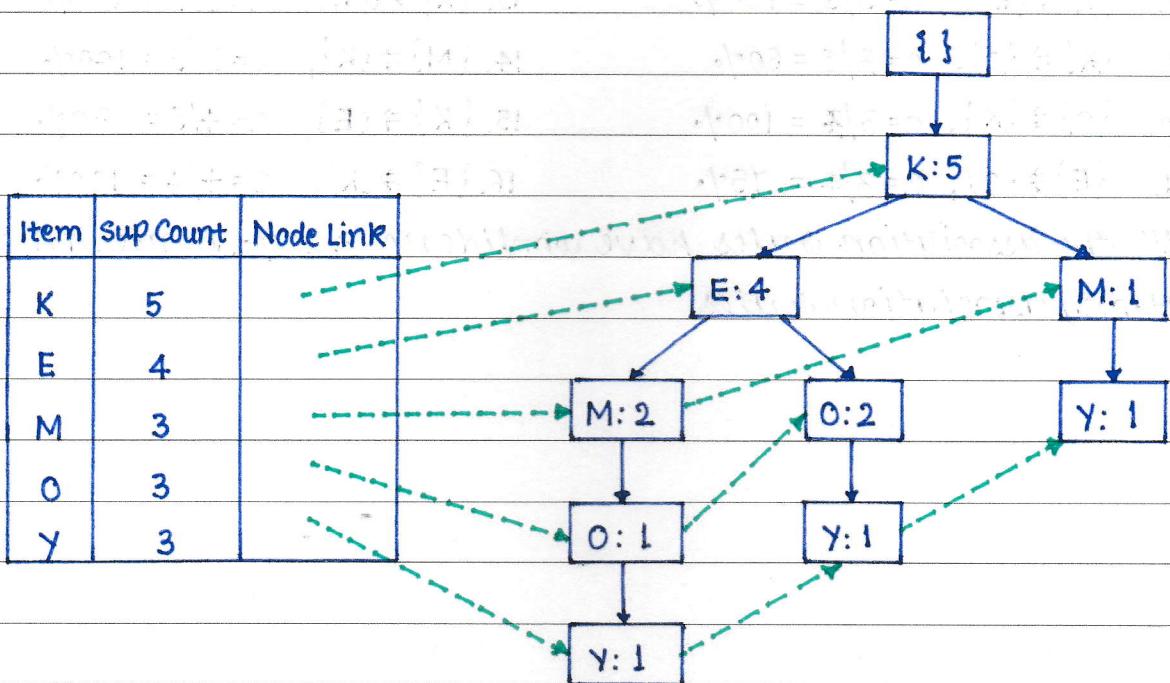
(ii) inserting {K,E,O,Y}



(iii) inserting {K,E,M}



(iv) Inserting  $\{K, M, Y\}$  and  $\{K, E, O\}$ , we have the final tree,



Therefore, conditional pattern base and conditional FP-tree are,

Item	conditional Pattern Base	Conditional FP-Tree
Y	$\{\{K, E, M, O:1\}, \{K, E, O:1\}, \{K, M:1\}\}$	$\langle K:3 \rangle$
O	$\{\{K, E, M:1\}, \{K, E:2\}\}$	$\langle K:3, E:3 \rangle$
M	$\{\{K, E:2\}, \{K:1\}\}$	$\langle K:3 \rangle$
E	$\{\{K:4\}\}$	$\langle K:4 \rangle$

so, frequent patterns generated,

Item	Frequent Patterns
Y	$\{K, Y:3\}$
O	$\{K, E, O:3\}, \{K, O:3\}, \{E, O:3\}$
M	$\{K, M:3\}$
E	$\{K, E:4\}$

Therefore, we have following association rules,

1.  $\{K, E\} \Rightarrow \{O\}$ ,  $c = 3/4 = \frac{75}{100}\%$
2.  $\{K, O\} \Rightarrow \{E\}$ ,  $c = 3/3 = 100\%$
3.  $\{E, O\} \Rightarrow \{K\}$ ,  $c = 3/3 = 100\%$
4.  $\{O\} \Rightarrow \{E, K\}$ ,  $c = 3/3 = 100\%$
5.  $\{E\} \Rightarrow \{O, K\}$ ,  $c = 3/4 = 75\%$
6.  $\{K\} \Rightarrow \{O, E\}$ ,  $c = 3/5 = 60\%$

- |   |   |
|---|---|
| 7. $\{K\} \Rightarrow \{Y\}$ , $C = 3/5 = 60\%$   | 12. $\{O\} \Rightarrow \{E\}$ , $C = 3/3 = 100\%$ |
| 8. $\{Y\} \Rightarrow \{K\}$ , $C = 3/3 = 100\%$  | 13. $\{K\} \Rightarrow \{M\}$ , $C = 3/5 = 60\%$  |
| 9. $\{K\} \Rightarrow \{O\}$ , $C = 3/5 = 60\%$   | 14. $\{M\} \Rightarrow \{K\}$ , $C = 3/3 = 100\%$ |
| 10. $\{O\} \Rightarrow \{K\}$ , $C = 3/3 = 100\%$ | 15. $\{K\} \Rightarrow \{E\}$ , $C = 4/5 = 80\%$  |
| 11. $\{E\} \Rightarrow \{O\}$ , $C = 3/4 = 75\%$  | 16. $\{E\} \Rightarrow \{K\}$ , $C = 4/4 = 100\%$ |

All the association rules have confidence  $\geq 20\%$ . Hence, all are strong association rules.

3. Why we use union in the formula of confidence instead of intersection?

The confidence of an association rule  $X \Rightarrow Y$  is given by:

$$\text{conf}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\text{supp}(X)}$$

The point to note here is that  $X$  and  $Y$  are itemsets not transactions. So,  $\text{supp}(X \cup Y)$  is the support where  $X$  and  $Y$  both itemsets occur. The union is used because here we are defining support in terms of sets of items and not events. We can use intersection if we define events  $E_x$  and  $E_y$ , where  $E_x$  denotes the event that the transaction contains itemset  $X$ . So, we can write,  
 $\text{support}(X \cup Y) = P(E_x \cap E_y)$ .

4. What is the physical significance of confidence and lift?

confidence of an association rule  $X \Rightarrow Y$  is given by,

$$\text{conf}(X \Rightarrow Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X)}$$

confidence, in simple words, is measure of the possibility that if an item  $X$  is bought, then the item  $Y$  will also be bought. It is conditional probability of  $X$  and  $Y$  both being purchased, given,  $X$  is already purchased.

Lift of an association rule  $X \Rightarrow Y$  is given by,

$$\text{lift}(X \Rightarrow Y) = \frac{\text{conf}(X \Rightarrow Y)}{\text{sup}(Y)} = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)\text{sup}(Y)}$$

Lift is a simple correlation measure, i.e. it signifies the dependency/independency of the two events, i.e. sale of X and sale of Y. In other words, it is a measure of how sale of X "lifts" the sale of Y. If  $\text{lift}(X \Rightarrow Y) < 1$ , it means the sale of X is going to decrease the sale of Y. If the value of lift is 1, it means, sale of Y is independent of sale of X. If value of lift  $> 1$ , it means, sale of X is likely going to increase the sale of Y.

### 5. What is the role of partitioning and sampling in apriori algorithm?

The partitioning and sampling algorithms are used to increase the efficiency of apriori algorithm. Basic apriori algorithm is slow as it requires complete database scan at each level of frequent itemset generation.

In partitioning the database is divided into n-non-overlapping partitions. The frequent itemsets local to each partition is found. Then all the frequent itemsets local to each partition is combined to form candidate itemset. Then the database is scanned for second time to find global frequent itemsets among those candidates.

In sampling, the random samples of database is picked and frequent itemsets are searched in s instead of complete database. But, since a random sample is searched, it is possible that some frequent itemsets are missed. In this way, we trade off some degree of accuracy against efficiency.