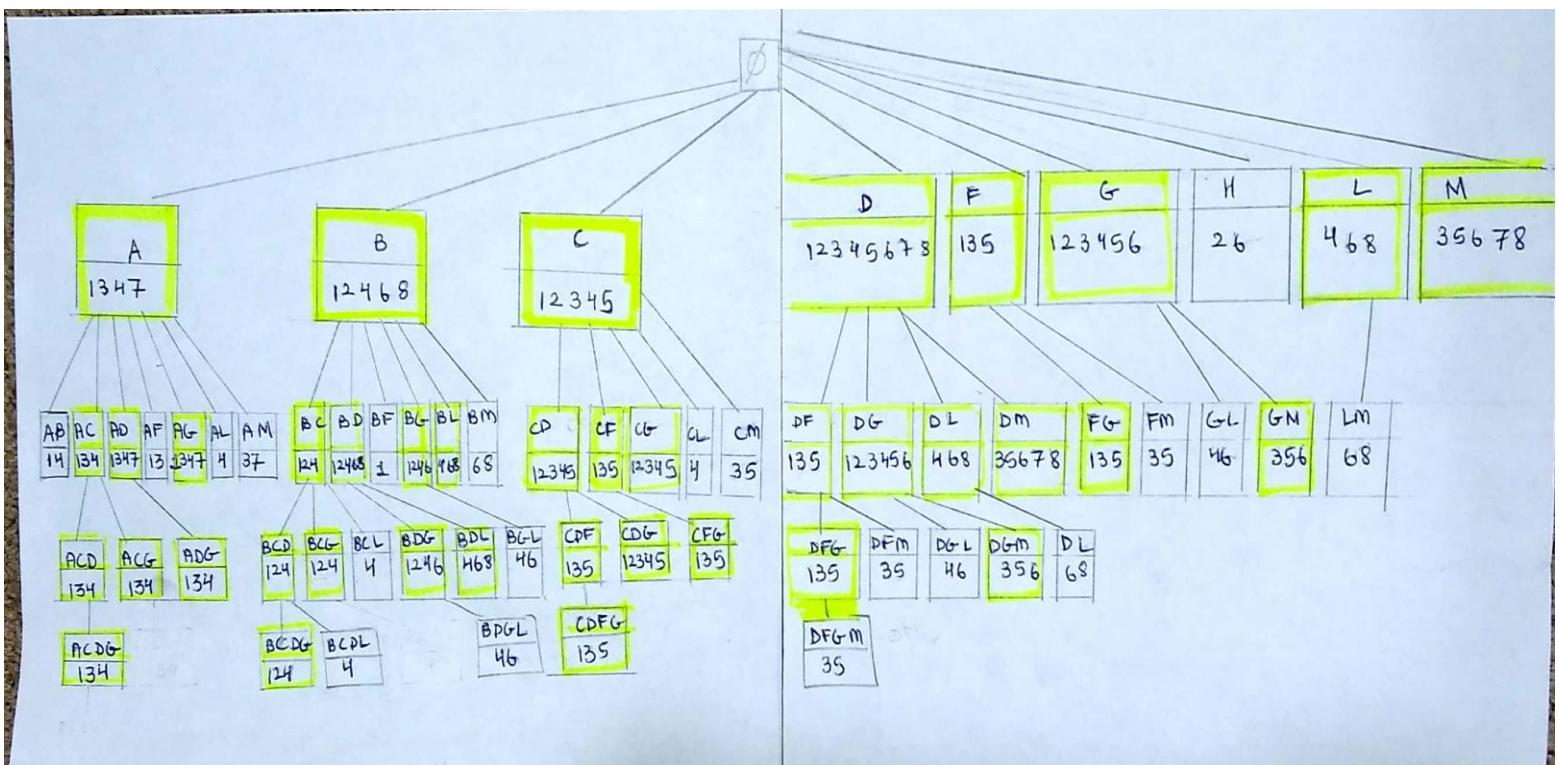


1. (Association Analysis) Consider the following transactions: (T1: A B C D F G), (T2: B C D G H), (T3: A C D F G K M), (T4: A B C D G L), (T5: C D F G M), (T6: B D G H L M), (T7: A D M), (T8: B D L M). Perform the following tasks using minimum support (minsup) value of 3.

- (15) Use Eclat algorithm to find all itemsets of size 3 or larger that are frequent. Show the execution of the algorithm and the results obtained.
- (15) Use GenMax algorithm to find all the maximal frequent itemsets in this dataset. Show the execution of the algorithm and the results obtained.
- (15) Use Charm algorithm to find all the closed itemsets. Show the execution of the algorithm and the results obtained.
- (15) Use FP-Tree algorithm to construct the FP Tree for the given set of transactions. Find all the frequent itemsets that contain the item "C". Show the final FP-Tree and main steps in constructing it. Also, show the execution of steps to find the frequent itemsets.

Answer 1)a)



I have dropped K as it's support is 1. Also drop H as its support is 2. At each step the itemsets with minimum support less than 3 are eliminated.

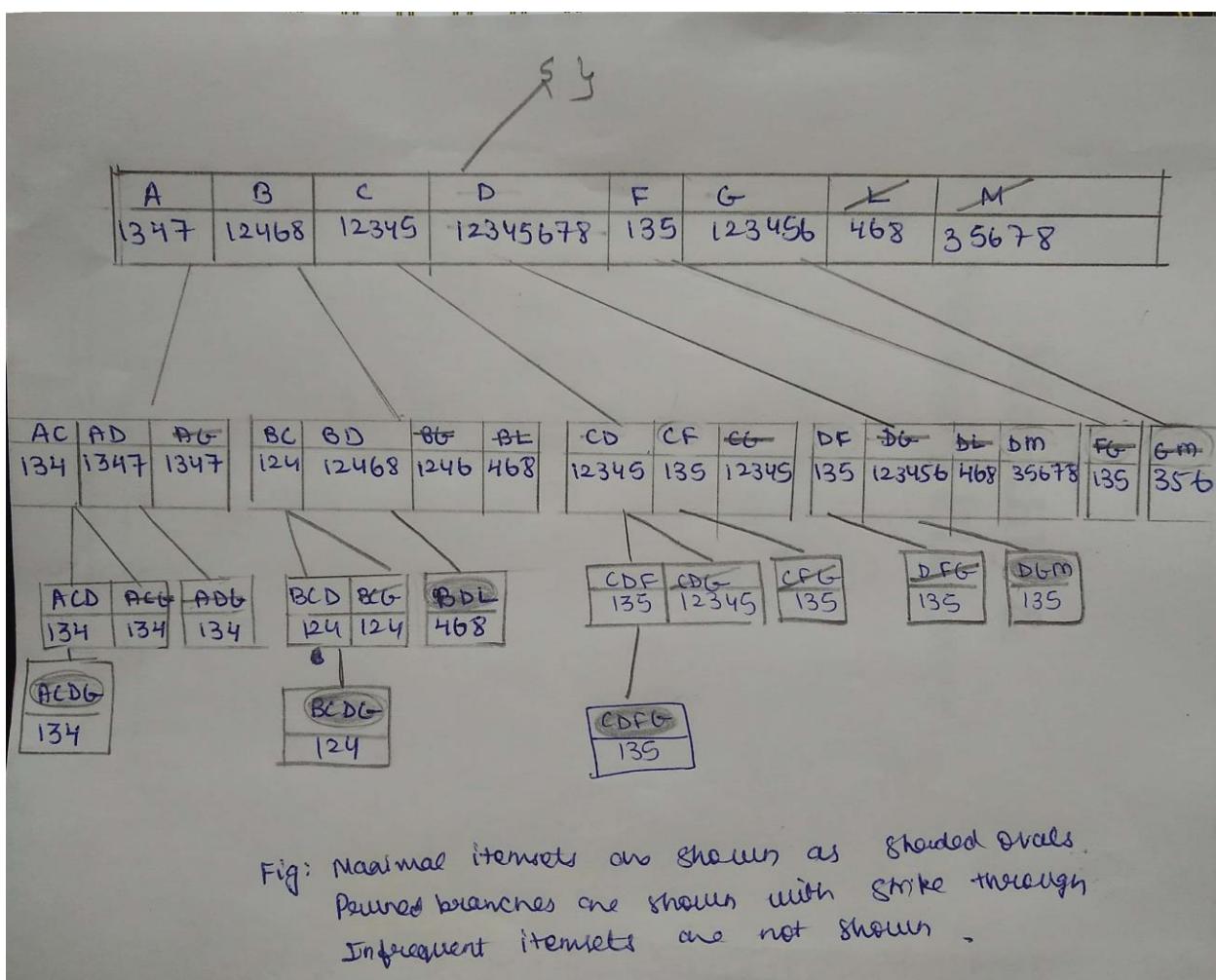
All those highlighted are frequent itemset

Results obtained-

The frequent itemsets are:

Minimum support	Itemsets
3	AC, BC, BL, CF, DF, DL, FG, GM, ACD, ACG, ADG, BCD, BCG, BDL, CDF, CFG, DFG, DGM ACDG, BCDG, CDFG, F, L
4	A, AD, AG, BG, BDG, BCDG
5	B, C, M, BD, CD, CG, CDG
6	G, DG
8	D

### Answer 1)b)



**The maximal frequent itemsets are:**

**ACDG, BCDG, BDL, CDFG, DGM**

Initially, all the itemsets with size 1 are written whose support is greater than 3. Then, itemsets of size 2 are generated and only the frequent itemsets are drawn in the graph. Following the **Depth First Strategy**. Itemset of size 3 and 4 are generated for the first branch.

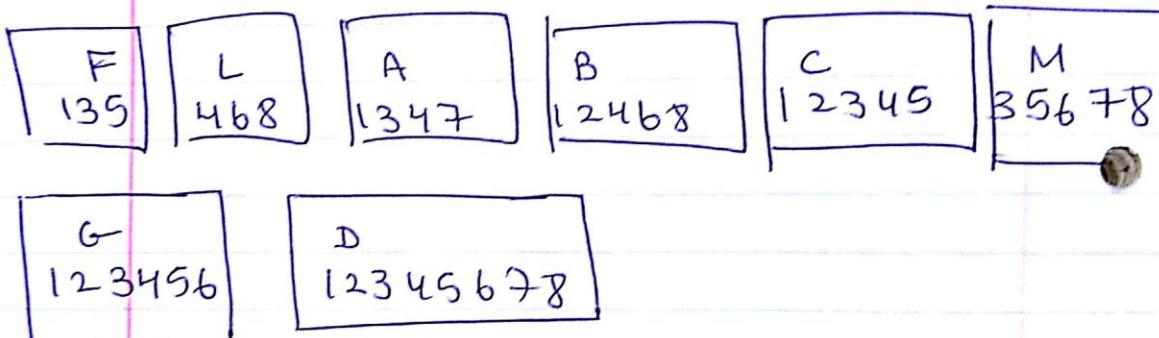
Now, the maximal frequent itemset is selected. Now, **backtracking** is done. All those itemsets are pruned which are the subset of this maximal frequent itemset.

These steps are repeated until we find all the maximal frequent itemsets.

Answer 1)c)

Initially the set of all closed itemsets,  $C$ , is empty.

Given any IT-pair set  $P = \{(x_i, t(x_i))\}$ , first sort them in the increasing order of their support



Now, the properties of cheen's algorithm are applied, we write only the frequent itemset,

As  $FL$ ,  $FA$ ,  $FB$ ,  $Fcm$  are infrequent, we ignore them ( $\text{min sup} < 3$ )

$$\text{Since } c(F) = c(FC)$$

$\therefore F$  can be replaced by  $FC$   
(property 1)

$$\text{Similarly } c(FC) = c(FCG)$$

$\therefore FC$  is replaced by  $FCG$

Similarly, the process is carried on when we get to A,  $c(A) \neq C(A)$ .

$t(AC) \subset t(A)$  is not true

$$\therefore t(AC) \neq t(A)$$

we apply the property 3 of Charnis Algorithm.

We add new extension  $P_A$ ,

If set  $P_A$  is not empty, we make recursive call to Charnis.

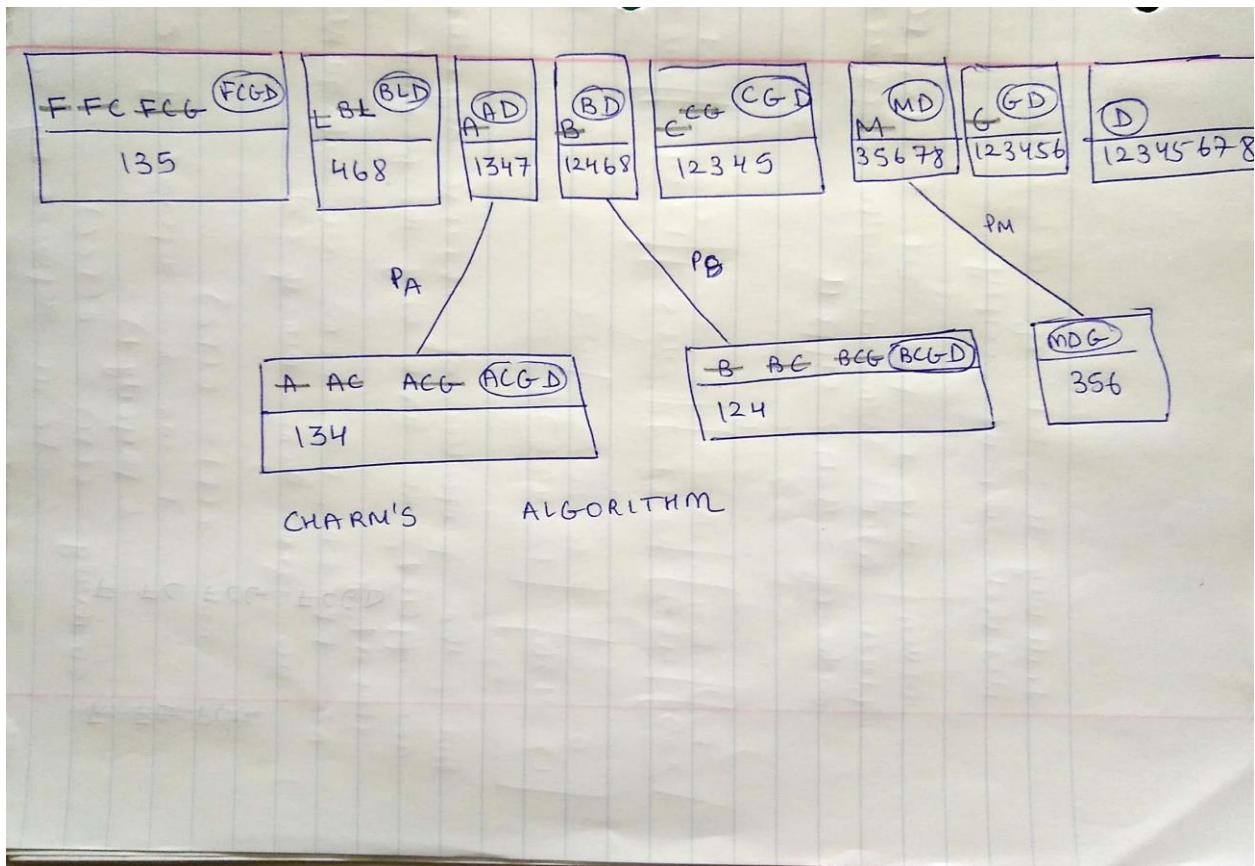
Hence, continue the same procedure.

Property (1) If  $t(X_i) = t(X_j)$ , then  $c(X_i) = c(X_j) = c(X_i \cup X_j)$ , which implies that we can replace every occurrence of  $X_i$  with  $X_i \cup X_j$  and prune the branch under  $X_j$  because its closure is identical to the closure of  $X_i \cup X_j$ .

Property (2) If  $t(X_i) \subset t(X_j)$ , then  $c(X_i) \neq c(X_j)$  but  $c(X_i) = c(X_i \cup X_j)$ , which means that we can replace every occurrence of  $X_i$  with  $X_i \cup X_j$ , but we cannot prune  $X_j$  because it generates a different closure. Note that if  $t(X_i) \supset t(X_j)$  then we simply interchange the role of  $X_i$  and  $X_j$ .

Property (3) If  $t(X_i) \neq t(X_j)$ , then  $c(X_i) \neq c(X_j) \neq c(X_i \cup X_j)$ . In this case we cannot remove either  $X_i$  or  $X_j$ , as each of them generates a different closure.

Thus, using the above properties of Charm's algorithms conclusions were drawn as follows.



The closed itemsets are:

**FCDG, BLD, AD, BD, CGD, MD, GD, D, ACGD, BCGD, MDG**

Answer 1)d)

## FP Tree Algorithm

①

### Step 1

Scan the database for the first time,  
find the frequent items (single item patterns)  
and order them into a set  $L$  in  
frequency descending order.

Given

Itemsets

$T_1$	A, B, C, D, F, G
$T_2$	B, C, D, G, H
$T_3$	A, C, D, F, G, K, M
$T_4$	A, B, C, D, G, L
$T_5$	C, D, F, G, M
$T_6$	B, D, G, H, L, M
$T_7$	A, D, M
$T_8$	B, D, L, M

### Performing Step 1:

$$L = \{ D: 8, G: 6, B: 5, C: 5, M: 5, A: 4, F: 3, \\ L: 3, H: 2, K: 1 \}$$

Now, drop the items which are not frequent  
i.e their support is less than min support = 3,

$$\therefore L = \{ D: 8, G: 6, B: 5, C: 5, M: 5, A: 4, F: 3, L: 3 \}$$

### Step 2

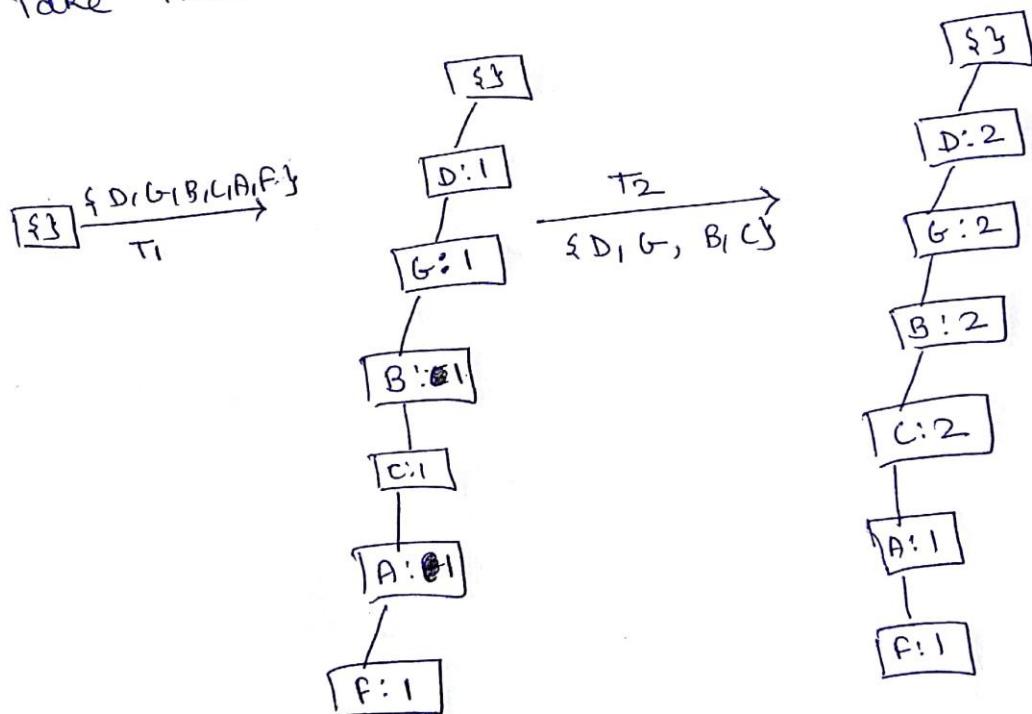
For each transaction, order its frequent items  
according to order in  $L$ .

Scan the database second time, construct  
FP tree by putting each frequency ordered  
transaction onto it.

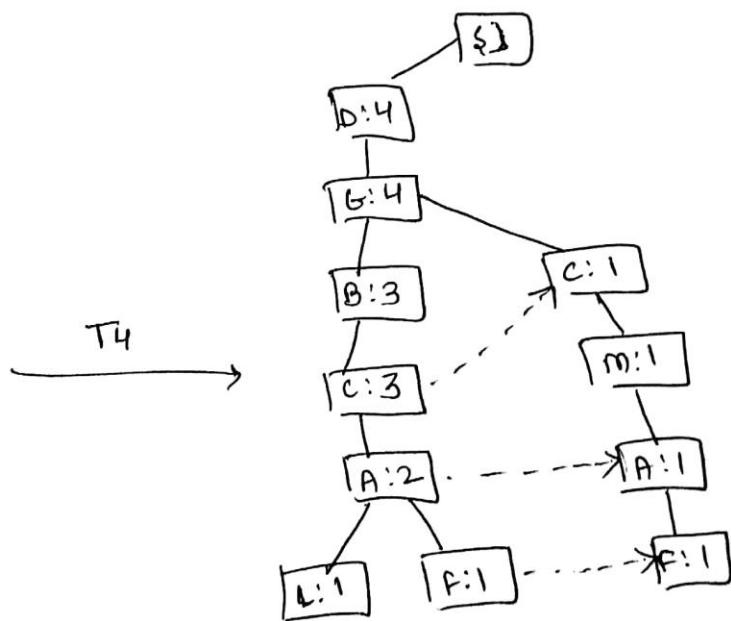
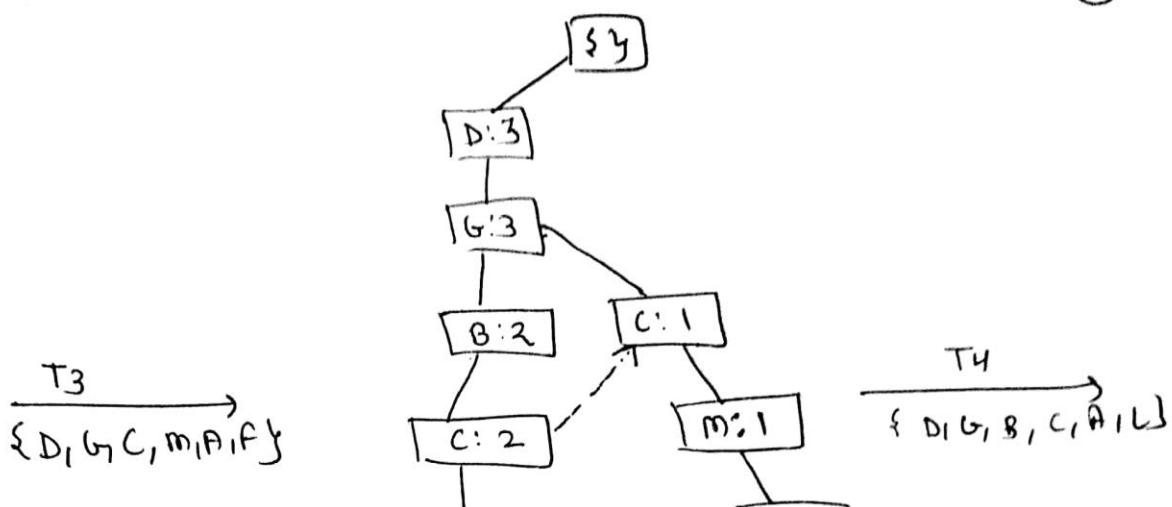
(2)

Transaction	Ordered frequent item
T <sub>1</sub>	D, G, B, C, A, F
T <sub>2</sub>	D, G, B, C
T <sub>3</sub>	D, G, C, M, A, F
T <sub>4</sub>	D, G, B, C, A, L
T <sub>5</sub>	D, G, C, M, F
T <sub>6</sub>	D, G, B, M, L
T <sub>7</sub>	D, M, A
T <sub>8</sub>	D, M, B, L

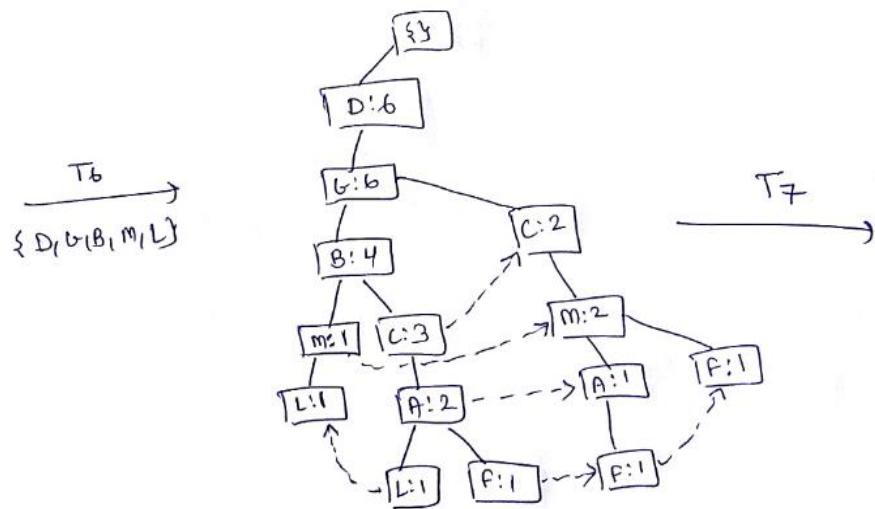
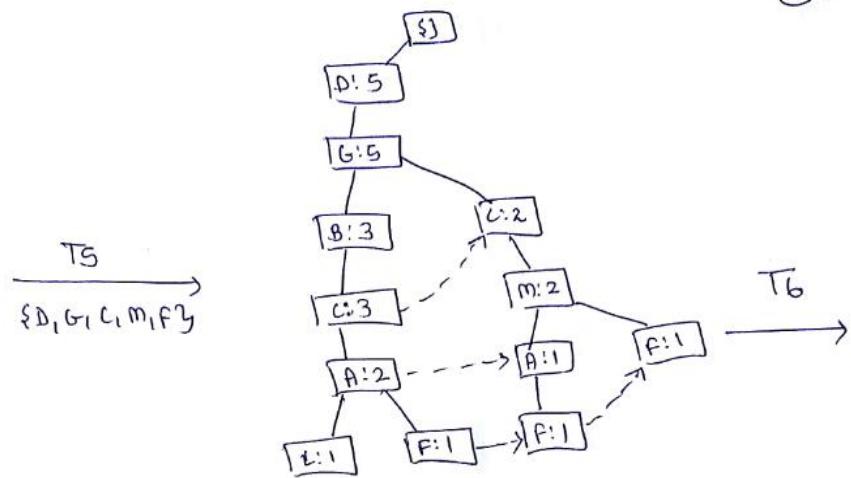
Now, construct the FP tree  
Take transaction 1 (T<sub>1</sub>)



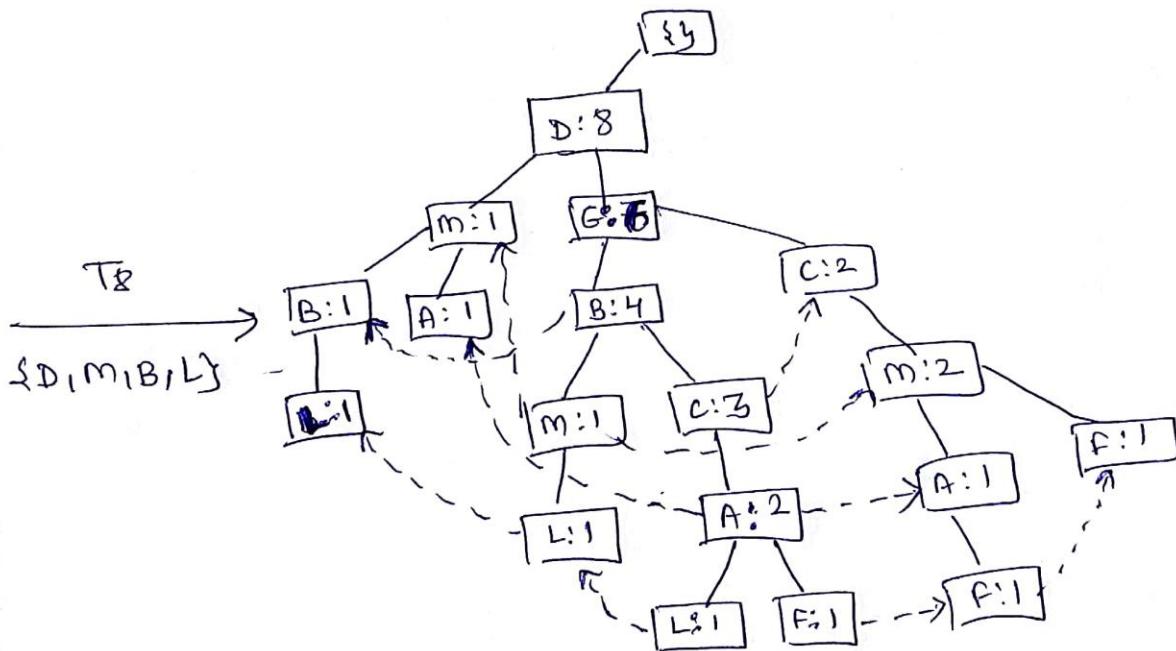
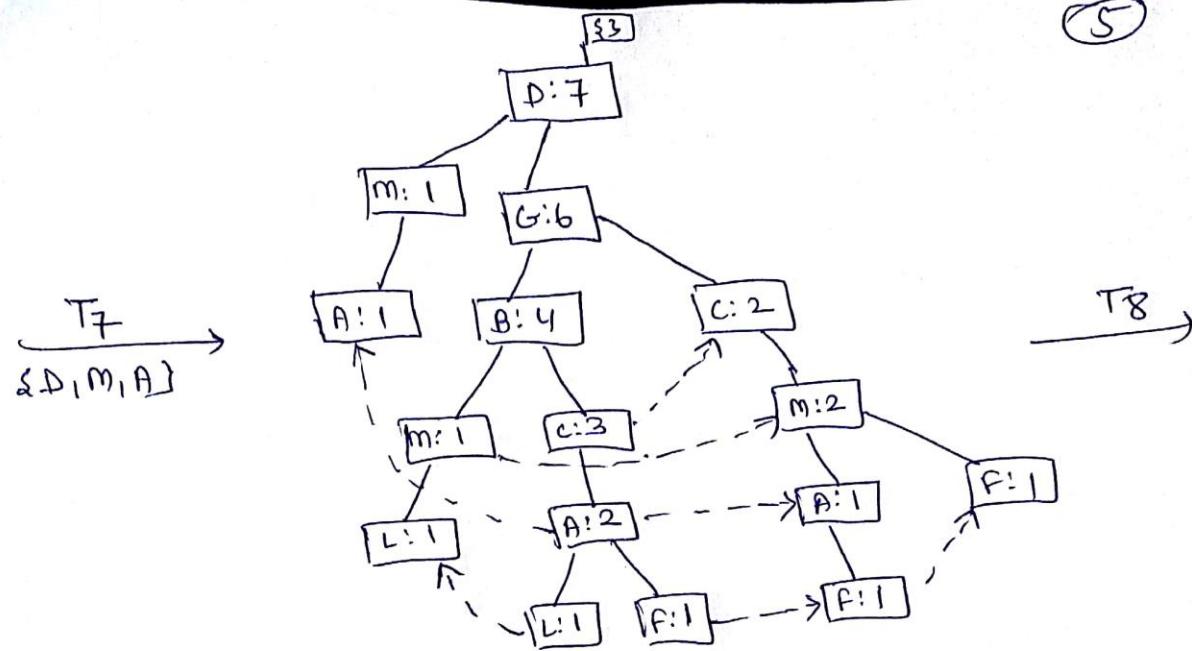
(3)



(4)

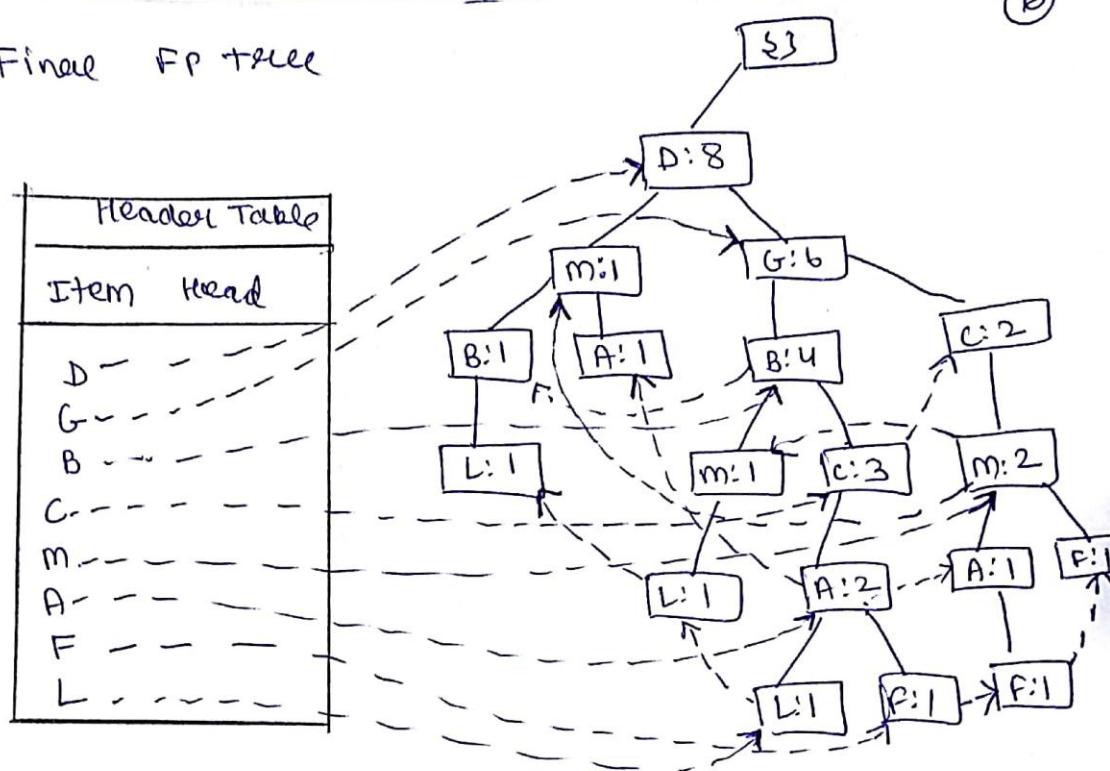


(5)



Final FP Tree

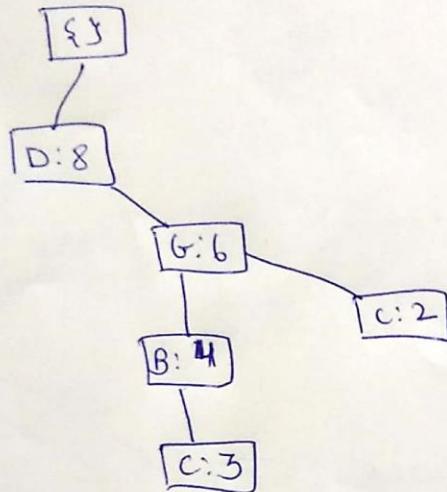
(b)



Now, the aim is to find out all the frequent itemset that contain C:

(7)

② Construct FP tree for the frequent item C



But since occurrence of  $C = 5 \therefore$ , the tree becomes

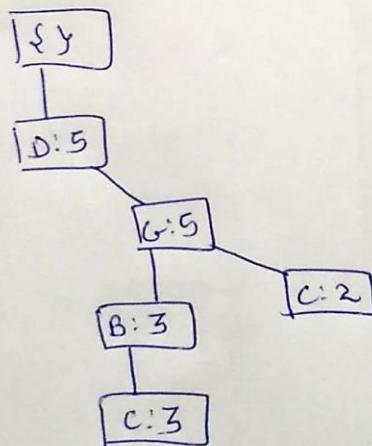
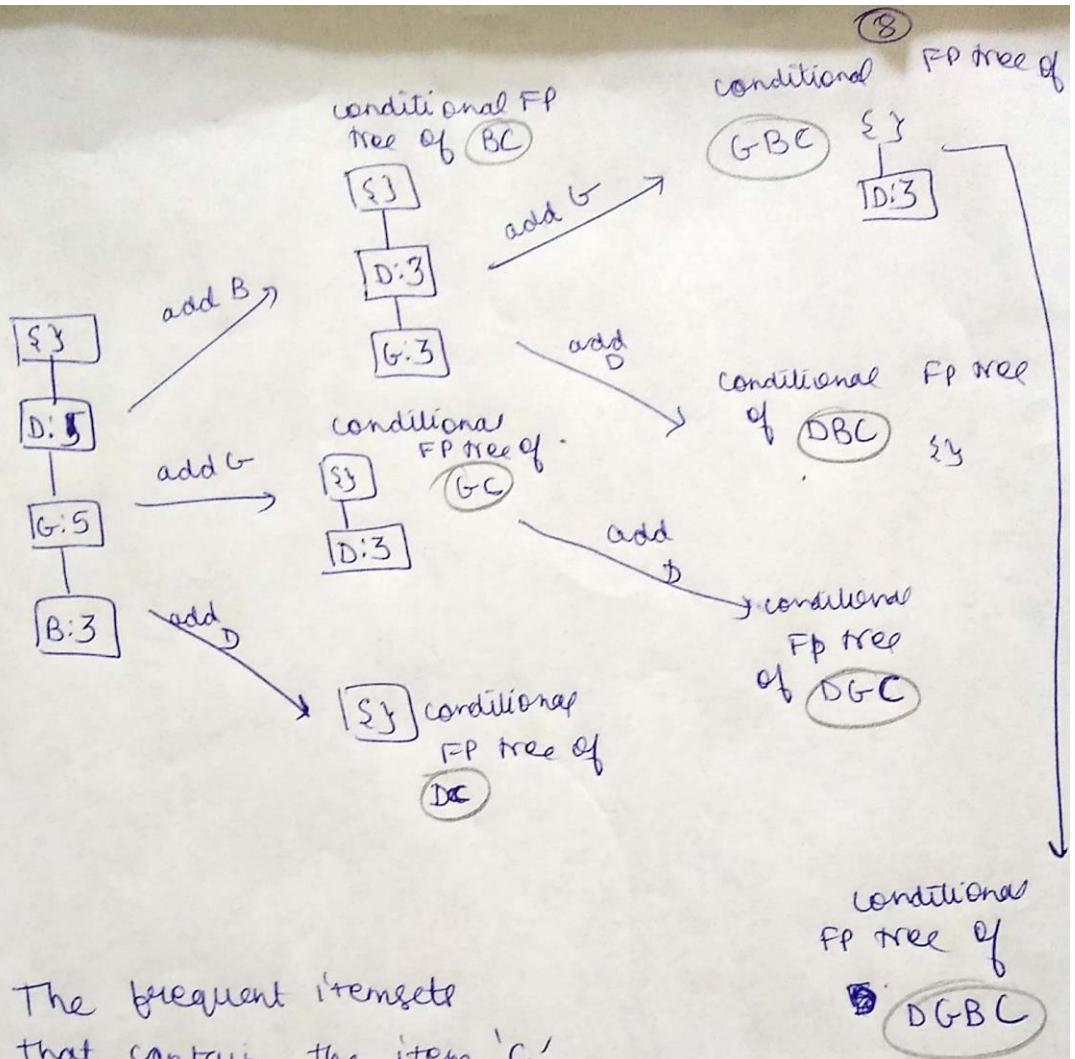


Fig: C - conditional FP tree



The frequent itemsets  
that contain the item 'C'

are:

BC, BCG, BCDG
GC, BCD
DC, CDG

2. Consider a dataset in which each row corresponds to a Plant and each column corresponds to a State in the US. Each entry in the table is marked as “1” when the Plant in the row is known to grow in the State for the column. An example table is shown below.
- (15) Without formally running any algorithm, and using only your intuition, select from this table two different closed itemsets and two different maximal frequent itemsets. Each of these itemsets must contain at least three items and also use  $\text{minsup}=3$ . Show the supporting sets of rows for each of these itemsets. Explain and show why each itemset has the property of being closed or maximal frequent.
  - (10) What interpretation can be given to the closed itemsets and what is a potential use for the knowledge of these closed itemsets?
  - (10) What interpretation can be given to the maximal frequent itemsets and what is a potential use for these maximal frequent itemsets?

	S1	S2	S3	S4	S5	S6
Pa	0	0	1	1	1	0
Pb	1	1	0	0	1	1
Pc	0	1	1	1	0	0
Pd	1	0	1	1	1	1
Pe	1	1	0	0	0	1
Pf	0	1	0	0	1	1
Pg	1	1	1	1	1	1
Ph	1	1	0	1	1	1
Pj	0	0	1	1	1	0

Answer 2)a)

①

Ans 5)

Firstly all the frequent itemsets  
are found:-

$$\begin{array}{l} \{S_1\} = 5 \\ \{S_2\} = 6 \\ \{S_3\} = 5 \\ \{S_4\} = 6 \\ \{S_5\} = 7 \\ \{S_6\} = 6 \end{array} \quad \left. \begin{array}{l} \text{all of them} \\ \text{have support} \\ \text{greater than the} \\ \text{min sup, } \therefore \text{they} \\ \text{are frequent.} \end{array} \right\}$$

Now, we find frequent itemsets of  
size 2.

itemset support

$$\{S_1, S_2\} = 4$$

$$\{S_1, S_3\} = 2 \quad [\text{Not frequent}]$$

$$\{S_1, S_4\} = 3$$

$$\{S_1, S_5\} = 4$$

$$\{S_1, S_6\} = 5$$

$$\{S_2, S_3\} = 2 \quad [\text{Not frequent}]$$

$$\{S_2, S_4\} = 3$$

$$\{S_2, S_5\} = 4$$

$$\{S_2, S_6\} = 5$$

$$\{S_3, S_4\} = 5$$

(2)

$$\{S_3, S_5\} = 4$$

$$\{S_3, S_6\} = 2 \quad [\text{Not frequent}]$$

$$\{S_4, S_5\} = 5$$

$$\{S_4, S_6\} = 3$$

$$\{S_5, S_6\} = 5$$

Now, we find frequent items of size 3:

$$\{S_1, S_2, S_3\} = 1 \quad [\text{Not frequent}]$$

$$\{S_1, S_2, S_4\} = 2 \quad [\text{Not frequent}]$$

$$\{S_1, S_2, S_5\} = 3$$

$$\{S_1, S_2, S_6\} = 4$$

$$\{S_1, S_3, S_4\} = 2 \quad [\text{Not frequent}]$$

$$\{S_1, S_3, S_5\} = 2 \quad [\text{Not frequent}]$$

$$\{S_1, S_3, S_6\} = 2 \quad [\text{Not frequent}]$$

$$\{S_1, S_4, S_5\} = 3$$

$$\{S_1, S_4, S_6\} = 3$$

$$\{S_1, S_5, S_6\} = 4$$

$$\{S_2, S_3, S_4\} = 2 \quad [\text{Not frequent}]$$

$$\{S_2, S_3, S_5\} = 1 \quad [\text{Not frequent}]$$

(3)

$$\{S_2, S_3, S_6\} = 1 \text{ [not frequent]}$$

$$\{S_2, S_4, S_5\} = 2 \text{ [not frequent]}$$

$$\{S_2, S_4, S_6\} = 2 \text{ [not frequent]}$$

$$\{S_2, S_4, S_5\} :$$

$$\{S_2, S_5, S_6\} = 3$$

$$\{S_3, S_4, S_5\} = 4$$

$$\{S_3, S_4, S_6\} = 2 \text{ [not frequent]}$$

$$\{S_3, S_5, S_6\} = 2 \text{ [not frequent]}$$

$$\{S_4, S_5, S_6\} = 3$$

Itemset which are frequent of size 4:

$$\{S_1, S_2, S_3, S_4\} = 1 \text{ [not frequent]}$$

$$\{S_1, S_2, S_3, S_5\} = 1 \text{ [not frequent]}$$

$$\{S_1, S_2, S_3, S_6\} = 1 \text{ [not frequent]}$$

$$\{S_1, S_2, S_4, S_5\} = 2 \text{ [not frequent]}$$

$$\{S_1, S_2, S_4, S_6\} = 2 \text{ [not frequent]}$$

$$\{S_1, S_2, S_5, S_6\} = 3$$

$$\{S_1, S_3, S_4, S_5\} = 2 \text{ [not frequent]}$$

$$\{S_1, S_3, S_4, S_6\} = 2 \text{ [not frequent]}$$

(n)

$$\{S_1, S_4, S_5, S_6\} = 3$$

$$\{S_2, S_3, S_4, S_5\} = 2 \text{ [Not frequent]}$$

$$\{S_2, S_3, S_4, S_6\} = 1 \text{ [Not frequent]}$$

$$\{S_2, S_4, S_5, S_6\} = 2 \text{ [Not frequent]}$$

$$\{S_3, S_4, S_5, S_6\} = 2 \text{ [Not frequent]}$$

$$\{S_1, S_2, S_3, S_4, S_5\} = 1 \text{ [Not frequent]}$$

$$\{S_1, S_2, S_3, S_4, S_6\} = 1 \text{ [Not frequent]}$$

$$\{S_1, S_2, S_3, S_5, S_6\} = 1 \text{ [Not frequent]}$$

$$\{S_1, S_2, S_4, S_5, S_6\} = 2 \text{ [Not frequent]}$$

$$\{S_1, S_3, S_4, S_5, S_6\} = 2 \text{ [Not frequent]}$$

$$\{S_2, S_3, S_4, S_5, S_6\} = 1 \text{ [Not frequent]}$$

$$\{S_1, S_2, S_3, S_4, S_5, S_6\} = 1 \text{ [Not frequent]}$$

(5)

### Frequent itemsets :

$$\{S_1\} = 5$$

$$\{S_3, S_4, S_5\} = 4$$

$$\{S_2\} = 6$$

$$\{S_4, S_5, S_6\} = 3$$

$$\{S_3\} = 5$$

$$\{S_1, S_2, S_5, S_6\} = 3$$

$$\{S_4\} = 6$$

$$\{S_1, S_4, S_5, S_6\} = 3$$

$$\{S_5\} = 7$$

$$\{S_6\} = 6$$

$$\{S_1, S_2\} = 4$$

$$\{S_1, S_4\} = 3$$

$$\{S_1, S_5\} = 4$$

$$\{S_1, S_6\} = 5$$

$$\{S_2, S_4\} = 3$$

$$\{S_2, S_5\} = 4$$

$$\{S_2, S_6\} = 5$$

$$\{S_3, S_4\} = 5$$

$$\{S_3, S_5\} = 4$$

$$\{S_4, S_5\} = 5$$

$$\{S_4, S_6\} = 3$$

$$\{S_5, S_6\} = 5$$

$$\{S_1, S_2, S_5\} = 3$$

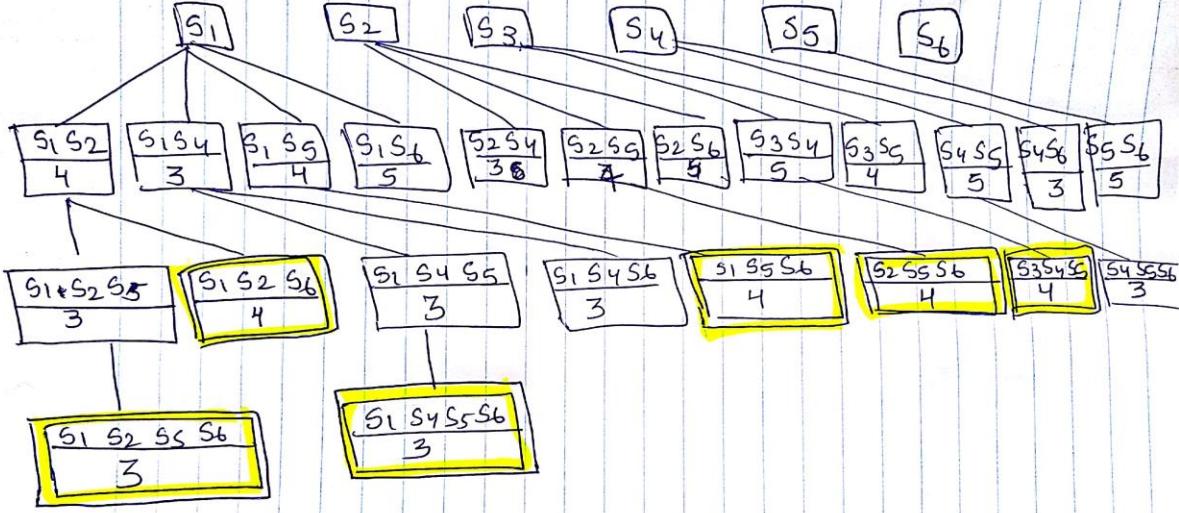
$$\{S_1, S_2, S_6\} = 4$$

$$\{S_1, S_4, S_5\} = 3$$

$$\{S_1, S_4, S_6\} = 3$$

$$\{S_1, S_5, S_6\} = 4$$

$$\{S_2, S_5, S_6\} = 3$$



The highlighted boxes contain maximal frequent itemsets. The maximal frequent itemsets are those which any of its supersets becomes infrequent.

$\{S_1, S_2, S_5, S_6\}$  is a maximal frequent itemset.

	$S_1$	$S_2$	$S_5$	$S_6$
Pb	1	1	1	1
Pg	1	1	1	1
Ph	1	1	1	1

Adding any other item such as  $S_3$

	$S_1$	$S_2$	$S_5$	$S_6$	$S_3$
Pb	1	1	1	1	0
Pg	1	1	1	1	0
Ph	1	1	1	1	0

The support goes down from 3 to 0,  
i.e. the superset  $\{S_1, S_2, S_5, S_6, S_3\}$   
becomes infrequent.

Adding any other item e.g.:  $S_4$

	$S_1$	$S_2$	$S_5$	$S_6$	$S_4$
Pb	1	1	1	1	0
Pg	1	1	1	1	1
Ph	1	1	1	1	1

The support goes down from 3 to 2,  
hence superset  $\{S_1, S_2, S_5, S_6, S_4\}$   
is infrequent.

Thus,  $\{S_1, S_2, S_5, S_6\}$  is a maximal frequent itemset.

Similarly

$\{S_1, S_4, S_5, S_6\}$  is a  
maximal frequent itemset

	$S_1$	$S_4$	$S_5$	$S_6$	
$P_d$	1	1	1	1	
$P_g$	1	1	1	1	
$P_h$	1	1	1	1	

Adding any other item  
suppose  $S_2$  is added

	$S_1$	$S_4$	$S_5$	$S_6$	$S_2$	
$P_d$	1	1	1	1	0	
$P_g$	1	1	1	1	1	
$P_h$	1	1	1	1	1	

Support goes down from 3 to 2,  
which is below the minsup threshold  
Hence the superset of  $\{S_1, S_4, S_5, S_6\}$   
which is  $\{S_2, S_1, S_4, S_5, S_6\}$  is  
infrequent -

Thus,  $\{S_1, S_4, S_5, S_6\}$  is  
maximal frequent.

Similarly other items are added and  
support definitely goes below the  
minsup.

$\{S_1 S_2 S_6\}$  is a closed itemset

	$S_1$	$S_2$	$S_6$
Pb	1	1	1
Pe	1	1	1
Pg	1	1	1
Ph	1	1	1

Adding another item for example  $S_5$  to it

	$S_1$	$S_2$	$S_6$	$S_5$
Pb	1	1	1	1
Pe	1	1	1	0
Pg	1	1	1	1
Ph	1	1	1	1

The support goes down from 4 to 3, but is still above the minsup threshold, i.e. the item is still frequent.

$\therefore \{S_1 S_2 S_6\}$  is a closed itemset as any of its superset has lower support than it.

$\{S_1, S_5, S_6\}$  is also a closed itemset.

	$S_1$	$S_5$	$S_6$
$P_b$	1	1	1
$P_d$	1	1	1
$P_g$	1	1	1
$P_h$	1	1	

Adding another item such as  $S_3$

	$S_1$	$S_5$	$S_6$	$S_3$
$P_b$	1	1	1	0
$P_d$	1	1	1	1
$P_g$	1	1	1	1
$P_h$	1	1	1	1

support goes down from 4 to 3,  
but does not go below the minsup.  
 $\therefore$  the  $\{S_1, S_5, S_6\}$  is a closed itemset.

As any of its superset has minimum support than it.

Similarly adding or  $S_2$  or  $S_4$  or  $S_4$  will also reduce min support.

**Answer 2)b)** An itemset is closed if none of its immediate supersets has the same support as the itemset

Use of knowledge of these closed itemsets:

- Having the closed itemsets we can **generate maximal frequent itemset as well as generate association rules.**
- Since, the datasets in real life are very large and finding all the frequent itemsets will just be resulting in **redundancy** of most of the sets. This redundancy of itemsets needs to be removed. This is made possible by the closed itemsets
- Knowledge of closed itemsets gives us the exact items we need to work on instead of giving the repeated itemsets.
- Another advantage of closed itemset is that the execution of algorithm to find the closed itemset (using CHARMS Algorithm) requires **less execution time** than to find the frequent itemset (using ECLAT algorithm)
- This difference in execution time is noticeable as we move towards larger dataset
- To generate rules, all we require it.

**Answer 2)c) An itemset is maximal frequent if none of its immediate supersets is frequent**

The potential use for the knowledge of frequent itemsets:

- Using The maximal frequent itemsets, we can find out frequent itemsets.
- Once we find maximal frequent itemset , we can generate all frequent itemset in **a single scan**. Because every subset of frequent itemset is frequent.
- This **saves time and space**, thus increases the efficiency
- It can be interpreted as a compact representation of all the frequent itemsets.
- There are only few maximal frequent itemsets thus they require the least space to store in the memory
- Another advantage of maximal frequent itemset is that the execution of algorithm to find the maximal frequent itemset (using GenMax Algorithm) requires **less execution time** than to find the frequent itemset (using ECLAT algorithm)



