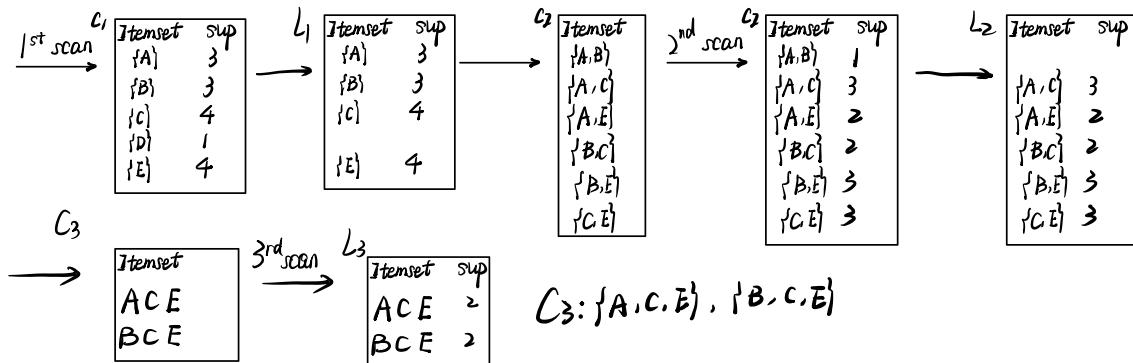


**Q1 [15 Marks]**

Given the transaction database below, set the minimum support count to 2 and the minimum confidence level to 60% to find the strong association rule. Generate the set  $C_3$  of the candidate 3-itemset, using pruning on Apriori principle.

| TID | Item    |
|-----|---------|
| T1  | A,C,D   |
| T2  | B,C,E   |
| T3  | A,B,C,E |
| T4  | B,E     |
| T5  | A,C,E   |

$sup_{min}=2$



step 1: Generate frequent 1-itemsets:  $\{A\}, \{B\}, \{C\}, \{E\}$

step 2: Generate frequent 2-itemsets:  $\{A, C\}, \{A, E\}, \{B, C\}, \{B, E\}, \{C, E\}$

step 3: Generate frequent 3-itemsets (apply pruning based on the Apriori principle)

$C_3: \{A, C, E\}, \{B, C, E\}$

step 4: calculate.

$$A \rightarrow C, A \rightarrow E, C \rightarrow E, C \rightarrow A, E \rightarrow A, E \rightarrow C, B \rightarrow C, B \rightarrow E, C \rightarrow B, E \rightarrow B$$

$$\frac{2}{3} \quad \frac{2}{3} \quad \frac{3}{4} \quad \frac{3}{4} \quad \frac{2}{7} \quad \frac{3}{7} \quad \frac{2}{3} \quad \frac{3}{3} \quad \frac{2}{4} \quad \frac{3}{4}$$

$$A \rightarrow C, E, C \rightarrow A, E, E \rightarrow A, C, C \rightarrow A, A, E \rightarrow C, A, C \rightarrow E$$

$$\frac{2}{3} \quad \frac{2}{4} \quad \frac{2}{4} \quad \frac{2}{3} \quad \frac{2}{2} \quad \frac{2}{2} \quad \frac{2}{3}$$

$$B \rightarrow C, E, C \rightarrow B, E, E \rightarrow B, C, C \rightarrow B, B \rightarrow C, B, C \rightarrow E.$$

$$\frac{2}{3} \quad \frac{2}{4} \quad \frac{2}{4} \quad \frac{2}{3} \quad \frac{2}{2} \quad \frac{2}{3} \quad \frac{2}{2}$$

confidence > 60%, the strong association rules are:

$$A \rightarrow C, A \rightarrow E, C \rightarrow E, C \rightarrow A, E \rightarrow C, B \rightarrow C, B \rightarrow E, E \rightarrow B$$

$$A \rightarrow C, E, C, E \rightarrow A, A, E \rightarrow C, A, C \rightarrow E, B \rightarrow C, E, C, E \rightarrow B, B, E \rightarrow C, B, C \rightarrow E.$$

## Q2 [15 Marks]

Reducing the transactions using dynamic hashing and pruning(DHP) algorithm. Set the minimum support count to 2.

Hash function bucket # =  $h(\{x\}) = ((\text{order of } x) * 10 + (\text{order of } y)) \% 7$

| TID | Item    |
|-----|---------|
| T1  | A,B,C   |
| T2  | B,D,E   |
| T3  | A,B,D,E |
| T4  | B,E     |

Items = A, B, C, D, E

Order = 1, 2, 3, 4, 5

The minimum support is 2

| Itemset | support |
|---------|---------|
| A       | 2       |
| B       | 4       |
| C       | 1       |
| D       | 2       |
| E       | 3       |

$C_1 \longrightarrow L_1$

Find all 2-itemset of each transaction

| TID | 2-itemset                     |
|-----|-------------------------------|
| T1  | {AB} {AC} {BC}                |
| T2  | {BD} {BE} {DE}                |
| T3  | {AB} {AD} {AE} {BD} {BE} {DE} |
| T4  | {BE}                          |

Hash function  $h(\{x\}) = ((\text{order of } x) * 10 + (\text{order of } y)) \% 7$

Hash table : map the 2-item of transaction into the hash table

|      |      |      |      |      |      |      |
|------|------|------|------|------|------|------|
| {AB} | {AE} | {BC} | {BD} | {BE} | {AB} | {AC} |
| {DE} |      |      | {BE} | {AB} |      |      |
| {BD} |      |      | {BE} |      |      |      |
| {DE} |      |      |      |      |      |      |

|          |   |   |   |   |   |   |
|----------|---|---|---|---|---|---|
| 1        | 1 | 1 | 4 | 3 | 2 | 1 |
| bucket 0 | 1 | 2 | 3 | 4 | 5 | 6 |

| $L_1 \times L_1$ | # in the bucket |
|------------------|-----------------|
| {AB}             | 2               |
| {AD}             | 1               |
| {AE}             | 1               |
| {BD}             | 4               |
| {BE}             | 3               |
| {DE}             | 4               |

|       |
|-------|
| $C_2$ |
| {AB}  |
| {BD}  |
| {BE}  |
| {DE}  |

| TID | Item    |
|-----|---------|
| T1  | A,B,C   |
| T2  | B,D,E   |
| T3  | A,B,D,E |
| T4  | B,E     |

$\xrightarrow{\quad}$

|      |           |
|------|-----------|
| {AB} | discard   |
| {BD} | keep {BD} |
| {BE} | keep {BE} |
| {DE} | discard   |

| $C_2$ | count | $L_2$ |
|-------|-------|-------|
| {AB}  | 2     | {AB}  |
| {BD}  | 2     | {BD}  |
| {BE}  | 2     | {BE}  |
| {DE}  | 2     | {DE}  |

$$D_3 = \langle T_2, BDE \rangle, \langle T_3, BDE \rangle$$

**Q3 [35 Marks]**

An itemset X is said to be a frequent itemset if the frequency count of X is at least a given support threshold.

An itemset Y is a **proper super-itemset** of X if  $X \subset Y$  and  $X \neq Y$ .

An itemset X is said to be a **closed frequent itemset** if (1) X is frequent and (2) there exists no proper super itemset Y of X such that Y is frequent and Y has the same frequency count as X.

An itemset X is said to be a **maximal frequent itemset** if (1) X is frequent and (2) there exists no proper super itemset Y of X such that Y is frequent.

Let  $F$  be the set of (traditional) frequent itemsets without specifying the frequency of itemsets.

Let  $F_c$  be the set of (traditional) frequent itemsets each of which is associated with

The following shows six transactions with four items. Each row corresponds to a transaction where 1 corresponds to a presence of an item and 0 corresponds to an absence.

|   | A | B | C | D |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 |
| 4 | 1 | 0 | 1 | 1 |
| 5 | 1 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 1 |

Suppose that the support threshold is 2.

(a) (i) What is  $F_c$ ? (ii) What is  $C_c$ ? (iii) What is  $M_c$ ? (5 Marks)

(b) (i) What are the advantages and the disadvantages of using closed frequent itemsets compared with traditional frequent itemsets? (5 Marks)

(ii) What are the advantages and the disadvantages of using closed frequent itemsets compared with maximal frequent itemsets? (5 Marks)

(c) Please adapt algorithm FP-growth with the use of the FP-tree to find all closed frequent itemset. Please write down how to adapt algorithm FP-growth and illustrate the adapted algorithm with the above example. (20 Marks)

a frequency in the dataset.

For example, if there are three frequent itemsets,  $\{I_1\}$  with frequency 4,  $\{I_2\}$  with frequency 5, and  $\{I_1, I_2\}$  with frequency 3,  $F = \{\{I_1\}, \{I_2\}, \{I_1, I_2\}\}$  and  $F_c = \langle \{I_1\}, 4 \rangle, \langle \{I_2\}, 5 \rangle, \langle \{I_1, I_2\}, 3 \rangle$ .

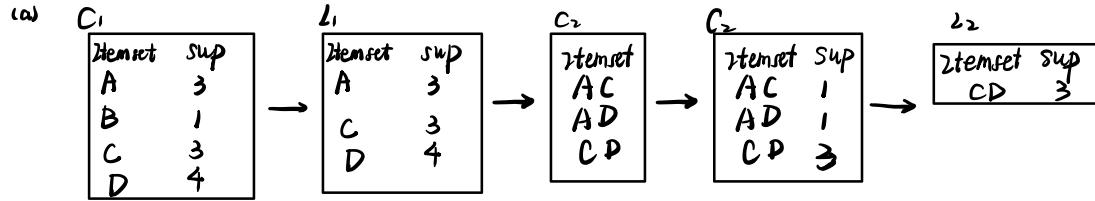
Similarly, let  $C$  be the set of closed frequent itemsets without specifying the frequency of itemsets.

Let  $C_c$  be the set of closed frequent itemsets each of which is associated with a frequency in the dataset.

Let  $M$  be the set of maximal frequent itemsets without specifying the frequency of itemsets.

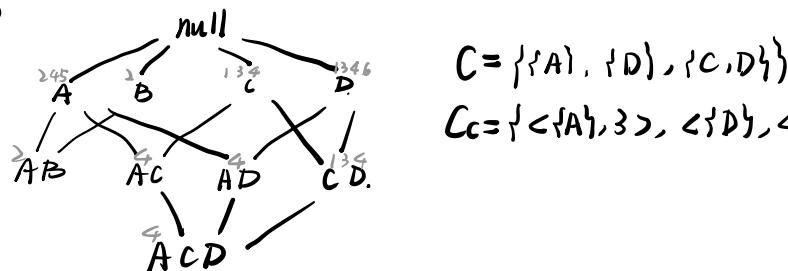
Let  $M_c$  be the set of maximal frequent itemsets each of which is associated with a frequency in the dataset.

| Tran | Items |
|------|-------|
| 1    | CD    |
| 2    | AB    |
| 3    | CD    |
| 4    | ACD   |
| 5    | A     |
| 6    | D     |



(ii)  $F_c = \langle \{A\}, 3 \rangle, \langle \{C\}, 3 \rangle, \langle \{D\}, 4 \rangle, \langle \{C, D\}, 3 \rangle$

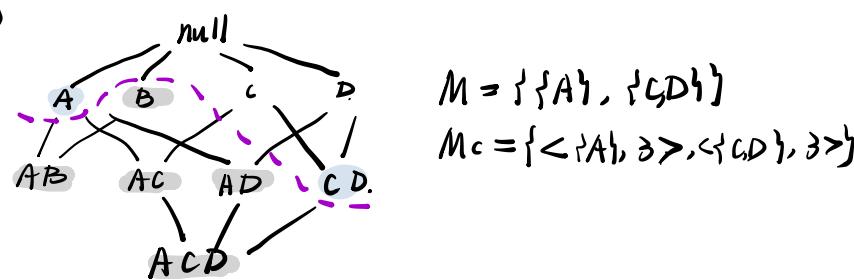
(iii)



$$C = \{\{A\}, \{D\}, \{C, D\}\}$$

$$C_c = \langle \{A\}, 3 \rangle, \langle \{D\}, 4 \rangle, \langle \{C, D\}, 3 \rangle$$

(iv)



$$M = \{\{A\}, \{CD\}\}$$

$$M_c = \langle \{A\}, 3 \rangle, \langle \{CD\}, 3 \rangle$$

(b) ii) Compared with traditional frequent itemsets

advantages: Closed frequent itemsets can reduce redundancy and computing is more efficient compared to traditional frequent itemsets, leading to faster results.

disadvantages: Using closed frequent itemsets may not capture all the information from the frequent itemsets and lose some useful frequent itemsets that are proper supersets and so on.

## ii) Compared with maximal frequent itemsets

advantages: Maximal frequent itemsets don't contain information about the support of their subsets and therefore an additional scan of transactions is required. On the other hand, closed frequent itemsets don't require this additional scan as they already capture the necessary support information.

disadvantages: Maximal frequent itemsets can inform the small set of itemsets that can derive all frequent itemsets but closed frequent itemsets can't. And implementation of closed frequent itemsets is relatively more complex.

And since maximal frequent itemset is also a part of closed itemset, the advantage and disadvantage of using closed frequent itemset instead of maximal frequent itemset is similar to the answer in (b)(i).

## (c) ① step 1: deduce the ordered frequent items

step 2: construct the FP-tree

step 3: construct the FP-conditional tree.

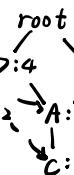
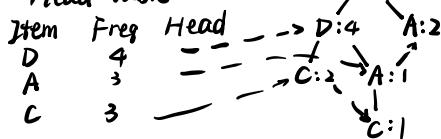
step 4: mine frequent patterns, identify closed itemsets for pruning and determine the frequent patterns.

② L<sub>1</sub>

| Itemset | Frequent |
|---------|----------|
| D       | 4        |
| A       | 3        |
| C       | 3        |

| Trans | Items | order frequent items |
|-------|-------|----------------------|
| 1     | CD    | DC                   |
| 2     | AB    | A                    |
| 3     | CD    | DC                   |
| 4     | ACD   | DAC                  |
| 5     | A     | A                    |
| 6     | D     | D                    |

③ Head table



cond FP-tree on C:3.

{(D:2, C:2), (D:1, C:1)} → {(D:2, C:1), (D:1, C:1)} → root

cond FP-tree on A:3

{(A:2), (D:1, A:1)} → {(A:3)} → root

cond FP-tree on D:4

{(D:4)} → root

| Item | cond pattern-base | Cond FP-tree. | Freq pattern. | closed Freq pattern |
|------|-------------------|---------------|---------------|---------------------|
| C    | {(D:2), (D, A:1)} | <D:3>         | {D, C:3}      | {D, C:3}            |
| A    | {(D:1)}           |               | {A:3}         | {A:3}               |
| D    |                   |               | {D:3}         | {D:3}               |

∴ {{C, D}, {A}, {D}}

**Q4 [35 Marks]**

A GSP Example: Suppose now we have 5 events: 'Upload Songs', 'Add Tags', 'Share', 'Listen' and 'Comment'. Let min-support be 40%. The sequence database of a Music Platform is shown in following table:

| Object | Sequence                                                     |
|--------|--------------------------------------------------------------|
| A      | <{'Upload Songs', 'Add Tags'}>                               |
| B      | <{'Upload Songs', 'Share'}>                                  |
| C      | <{'Upload Songs'}, {'Share', 'Listen'}>                      |
| D      | <{'Upload Songs'}, {'Upload Songs', 'Add Tags'}, {'Listen'}> |
| E      | <{'Listen'}, {'Add Tags', 'Comment'}, {'Share', 'Listen'}>   |

Please answer the following questions:

- (a) Make the first pass over the sequence database to yield all the 1-element frequent sequences and what is the corresponding support? (5 Marks)
- (b) Based on (a), do the 2-sequences Candidate Generation and Candidate Pruning.

(10 Marks)

- (c) What is the frequent 2-sequences based on the results of (b)? (5 Marks)
- (d) Based on (c), do the 3-sequences Candidate Generation and Candidate Pruning. When a sequence should be pruned, you need to explain why. (10 Marks)
- (e) What is the frequent 3-sequences based on the results of (d)? Please calculate the support. (5 Marks)

**Remember:** For frequent k-sequences, the support  $\geq$  min-support

**(a) Candidate 1-sequences are:**

<{'Upload Songs'}>, <{'Add Tags'}>, <{'Share'}>, <{'Listen'}>, <{'Comment'}>  
80% 60% 60% 60% 20%

$\therefore$  all the 1-element frequent sequences are: <{'Upload Songs'}>, <{'Add Tags'}>, <{'Share'}>, <{'Listen'}>.

**(b) Candidate 2-sequences are:**

<{'Upload Songs', 'Add Tags'}>, <{'Upload Songs', 'Share'}>, <{'Upload Songs', 'Listen'}>,  
 <{'Add Tags', 'Share'}>, <{'Add Tags', 'Listen'}>, <{'Share', 'Listen'}>,  
 <{'Upload Songs'}, {'Upload Songs'}>, <{'Upload Songs'}, {'Add Tags'}>,  
 <{'Upload Songs'}, {'Share'}>, <{'Upload Songs'}, {'Listen'}>.  
 <{'Add Tags'}, {'Upload Songs'}>, <{'Add Tags'}, {'Add Tags'}>,  
 <{'Add Tags'}, {'Share'}>, <{'Add Tags'}, {'Listen'}>,  
 <{'Share'}, {'Upload Songs'}>, <{'Share'}, {'Add Tags'}>,  
 <{'Share'}, {'Share'}>, <{'Share'}, {'Listen'}>,  
 <{'Listen'}, {'Upload Songs'}>, <{'Listen'}, {'Add Tags'}>,  
 <{'Listen'}, {'Share'}>, <{'Listen'}, {'Listen'}>, 22 sequences

After candidate pruning,

<{'Upload Songs', 'Add Tags'}>, <{'Upload Songs', 'Share'}>, <{'Upload Songs', 'Listen'}>,  
 <{'Add Tags', 'Share'}>, <{'Add Tags', 'Listen'}>, <{'Share', 'Listen'}>,  
 <{'Upload Songs'}, {'Upload Songs'}>, <{'Upload Songs'}, {'Add Tags'}>,  
 <{'Upload Songs'}, {'Share'}>, <{'Upload Songs'}, {'Listen'}>.  
 <{'Add Tags'}, {'Upload Songs'}>, <{'Add Tags'}, {'Add Tags'}>,  
 <{'Add Tags'}, {'Share'}>, <{'Add Tags'}, {'Listen'}>,  
 <{'Share'}, {'Upload Songs'}>, <{'Share'}, {'Add Tags'}>,  
 <{'Share'}, {'Share'}>, <{'Share'}, {'Listen'}>,  
 <{'Listen'}, {'Upload Songs'}>, <{'Listen'}, {'Add Tags'}>,  
 <{'Listen'}, {'Share'}>, <{'Listen'}, {'Listen'}>,

**(c) After candidate elimination, the remaining frequent 2-sequences are:**

<{'Upload Songs', 'Add Tags'}> (support=0.4)  
 <{'Share'}, {'Listen'}> (support=0.4)  
 <{'Upload Songs'}, {'Listen'}> (support=0.4)  
 <{'Add Tags'}, {'Listen'}> (support=0.4)

(d) Candidate 3-sequences :

$\langle \{\text{Upload Songs}, \text{Add Tags}\}, \{\text{listen}\} \rangle$ .

generated from  $\langle \{\text{Upload Songs}, \text{Add Tags}\} \rangle$  and  $\langle \{\text{Add Tags}\}, \{\text{listen}\} \rangle$

and should not be pruned because all 2-subsequences

$\langle \{\text{Upload Songs}\}, \{\text{listen}\} \rangle, \langle \{\text{Add Tags}\}, \{\text{listen}\} \rangle$  are frequent

$\langle \{\text{Upload Songs}\}, \{\text{share}, \text{listen}\} \rangle$

generated from  $\langle \{\text{Upload Songs}\}, \{\text{listen}\} \rangle$  and  $\langle \{\text{share}\}, \{\text{listen}\} \rangle$

and should be pruned because  $\langle \{\text{Upload Songs}\}, \{\text{share}\} \rangle$  is not frequent.

$\langle \{\text{Add Tags}\}, \{\text{share}, \text{listen}\} \rangle$

generated from  $\langle \{\text{Add Tags}\}, \{\text{listen}\} \rangle$  and  $\langle \{\text{share}\}, \{\text{listen}\} \rangle$

and should be pruned because  $\langle \{\text{Add Tags}\}, \{\text{share}\} \rangle$  is not frequent.

So, after pruning, the remaining 3-sequence:

$\langle \{\text{Upload Songs}, \text{Add Tags}\}, \{\text{listen}\} \rangle$

(e)  $\langle \{\text{Upload Songs}, \text{Add Tags}\}, \{\text{listen}\} \rangle$ . support = 0.2 < 0.4,

should be eliminated.

Thus, there are no 3-sequences left.