

## Homework 1

**50015627 JIANG Zhuoyang**

**Q1 [15 Marks]**

5/15

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

With Apriori, get the frequent k-itemsets and generate the candidate  $k+1$ -itemsets as below:

$C_1$		$L_1$	$C_2$	
Itemset	sup	prune	Itemset	sup
{A}	3	→	{A}	3
{B}	3		{B}	3
{C}	4		{C}	4
{D}	1		{E}	4
{E}	4			

2nd Scan →  $C_2$

$C_2$		$L_2$	$C_3$	
Itemset	sup	prune	Itemset	sup
{A, B}	1	→	{A, C}	3
{A, E}	3		{A, E}	2
{A, C}	2		{B, C}	2
{B, C}	2		{B, E}	3
{B, E}	3		{C, E}	3
{C, E}	3			

$(C_3 = \{I \mid I \in L_2 \text{ and } I \text{ with } AB \notin I\})$

3rd Scan →  $C_3$

$C_3$		$\Rightarrow C_3 = \{\{A, C, E\}, \{B, C, E\}\}$
Itemset	sup	
{A, C, E}	2	
{B, C, E}	2	

## Q2 [15 Marks]

15/15

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

C <sub>1</sub>	L <sub>1</sub>																						
<table border="1"> <thead> <tr> <th>Itemset</th> <th>sup</th> </tr> </thead> <tbody> <tr> <td>{A}</td> <td>2</td> </tr> <tr> <td>{B}</td> <td>4</td> </tr> <tr> <td>{C}</td> <td>1</td> </tr> <tr> <td>{D}</td> <td>2</td> </tr> <tr> <td>{E}</td> <td>3</td> </tr> </tbody> </table>	Itemset	sup	{A}	2	{B}	4	{C}	1	{D}	2	{E}	3	<table border="1"> <thead> <tr> <th>Itemset</th> <th>sup</th> </tr> </thead> <tbody> <tr> <td>{A}</td> <td>2</td> </tr> <tr> <td>{B}</td> <td>4</td> </tr> <tr> <td>{D}</td> <td>2</td> </tr> <tr> <td>{E}</td> <td>3</td> </tr> </tbody> </table>	Itemset	sup	{A}	2	{B}	4	{D}	2	{E}	3
Itemset	sup																						
{A}	2																						
{B}	4																						
{C}	1																						
{D}	2																						
{E}	3																						
Itemset	sup																						
{A}	2																						
{B}	4																						
{D}	2																						
{E}	3																						

Solution:  $L_1 = \{\{A\}, \{B\}, \{D\}, \{E\}\}$ , for all 2-itemsets in transactions  
creat the hash table below:

bucket	0	1	2	3	4	5	6
count	1	1	1	4	3	2	1
itemset	{A,D}	{A,E}	{B,C}	{B,D}	{B,E}	{A,B}	{A,C}

$\Rightarrow$  Generate  $C_2$ :

L <sub>1</sub> , N <sub>1</sub>	Itemset	#in bucket of the Itemset	prune	C <sub>2</sub>	$L_2$
	{A,B}	2		$C_2 = \{\{B,D\}, \{D,E\}, \{B,E\}, \{A,B\}\}$	
	{A,D}	1			
	{A,E}	1			
	{B,D}	4			
	{B,E}	3			
	{D,E}	4			

so that we have:

D<sub>2</sub>:

T<sub>1</sub>: {A, B}

$\Rightarrow$  Discarded

T<sub>2</sub>: {B, D}, {B,E}, {D,E}

$\Rightarrow$  keep {B, D, E}

(item A occurs just in one 2-itemset {A,B}. so that 3-itemset with A will be discarded.)

T<sub>3</sub>: {A, B}, {B, D}, {B, E}, {D, E}

$\Rightarrow$  keep {B, D, E}

T<sub>4</sub>: {B, E}

$\Rightarrow$  Discarded

$$\Rightarrow D_3 = \left\{ \{T_2, \{B, D, E\}\}, \{T_3, \{B, D, E\}\} \right\}$$

### Q3 [35 Marks]

- 1 An itemset X is said to be a frequent itemset if the frequency count of X is at least a given support threshold.
- 2 An itemset Y is a proper super-itemset of X if  $X \subset Y$  and  $X \neq Y$ .
- 3 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.
- 4 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.
- 5 Let  $F$  be the set of (traditional) frequent itemsets without specifying the frequency of itemsets.
- 6 Let  $F_c$  be the set of (traditional) frequent itemsets each of which is associated with 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  $F_c = \{\langle\{I_1\}, 4\rangle, \langle\{I_2\}, 5\rangle, \langle\{I_1, I_2\}, 3\rangle\}$ .
- 7 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.
- 8 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.

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
0	0	1	1
1	1	0	0
0	0	1	1
1	0	1	1
1	0	0	0
0	0	0	1

Suppose that the support threshold is 2.

Solution :

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

- i)  $F_c = \{\langle \{A\}, 3 \rangle, \langle \{C\}, 1 \rangle, \langle \{D\}, 4 \rangle, \langle \{C,D\}, 3 \rangle\}$  3/5  
ii)  $C_c = \{\langle \{A\}, 3 \rangle, \langle \{D\}, 4 \rangle, \langle \{C,D\}, 3 \rangle\}$   
iii)  $M_c = \{\langle \{A\}, 3 \rangle, \langle \{C,D\}, 3 \rangle\}$  no process

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

- Advantages: ① They provide a relatively compact representation of frequent itemsets. ② They contain the complete support information of their subsets. So that we can not only obtain the non-frequent itemsets but also obtain their support count without any additional pass.  
Disadvantages: Their representation is not as compact as Maximal frequent itemsets.

(ii) What are the advantages and the disadvantages of using closed frequent itemsets compared with maximal frequent itemsets? (5 Marks) 4/5

- Advantages: They can reach the most compact representation of frequent itemset in which all of the frequent itemsets can be obtained.  
Disadvantages: They don't contain the complete support information of their subset, and an additional pass is needed.

(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) 10/20

i) Construct the FP-Tree:

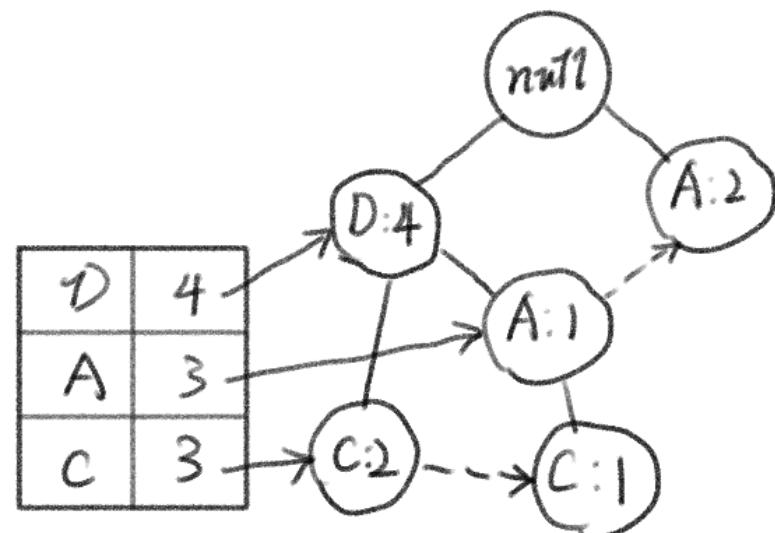
① Scan and generate  $L_1$ :  
 $\{D:4, A:3, C:3\}$

③ Scan and build the FP-Tree:

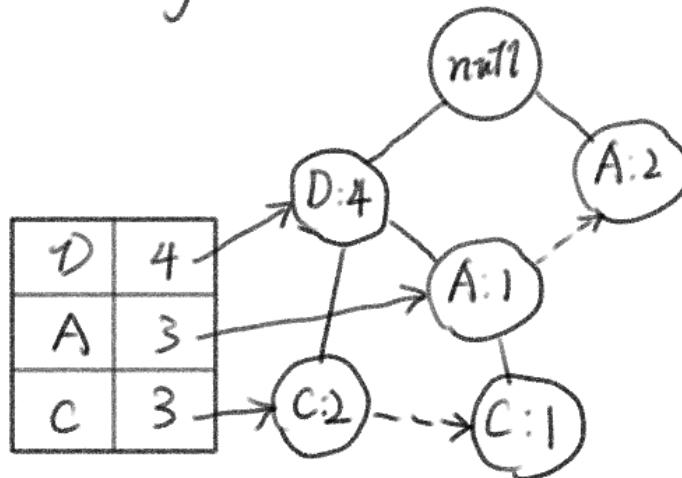
② Create the root node:

D	4
A	3
C	3

null



ii) FP-Growth for pattern mining:



① Start from the least frequent 1-pattern. Here we have C.

$\Rightarrow$  for C:

1° construct its conditional pattern base, for C, we have:

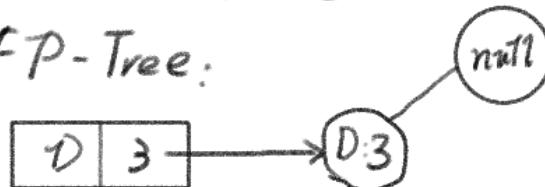
$\langle D, C:2 \rangle, \langle D, A, C:1 \rangle$

$\Rightarrow$  conditional pattern base:

$\langle D:2 \rangle, \langle D, A:1 \rangle$

2° Construct its conditional

FP-Tree:



3° Generate all frequent patterns for C:  $\{\langle D, C:3 \rangle\}$

② for A:  $\langle D, A:1 \rangle, \langle A:2 \rangle$

1° construct its conditional pattern base: null.

2° construct its conditional FP-Tree: null

③ for D:  $\langle D:4 \rangle$

1° construct its conditional pattern base: null.

2° construct its conditional FP-Tree: null

iii) So that, the resulted frequent itemsets:

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

$\Rightarrow$  the resulted closed frequent itemsets:

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

Q4 [35 Marks]

35/35 good

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: (2/5)

Object	Sequence
A	$\langle \{ \text{'Upload Songs'}, \text{'Add Tags'} \} \rangle$
B	$\langle \{ \text{'Upload Songs'}, \text{'Share'} \} \rangle$
C	$\langle \{ \text{'Upload Songs'} \}, \{ \text{'Share'}, \text{'Listen'} \} \rangle$
D	$\langle \{ \text{'Upload Songs'} \}, \{ \text{'Upload Songs'}, \text{'Add Tags'}, \text{'Listen'} \} \rangle$
E	$\langle \{ \text{'Listen'} \}, \{ \text{'Add Tags'}, \text{'Comment'} \}, \{ \text{'Share'}, \text{'Listen'} \} \rangle$

Use the first letter to represent each event: 'Upload Songs' = U; 'Add Tags' = A; 'Share' = S; 'Listen' = L; 'Comment' = C

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)

1-element frequent sequence & their support:

$\langle \{U\} \rangle : s = 80\%$ ;  $\langle \{A\} \rangle : s = 60\%$ ;  $\langle \{S\} \rangle : s = 60\%$ ;  $\langle \{L\} \rangle : s = 60\%$

(b) Based on (a), do the 2-sequences Candidate Generation and Candidate Pruning.

1) 2-element frequent sequence generation: (10 Marks)

$\langle \{U, A\} \rangle, \langle \{U, S\} \rangle, \langle \{U, L\} \rangle, \langle \{A, S\} \rangle, \langle \{A, L\} \rangle, \langle \{S, L\} \rangle,$   
 $\langle \{U\}, \{U\} \rangle, \langle \{U\}, \{A\} \rangle, \langle \{U\}, \{S\} \rangle, \langle \{U\}, \{L\} \rangle, \langle \{S\}, \{S\} \rangle, \langle \{S\}, \{U\} \rangle, \langle \{S\}, \{A\} \rangle, \langle \{S\}, \{L\} \rangle,$   
 $\langle \{A\}, \{A\} \rangle, \langle \{A\}, \{U\} \rangle, \langle \{A\}, \{S\} \rangle, \langle \{A\}, \{L\} \rangle, \langle \{L\}, \{L\} \rangle, \langle \{L\}, \{U\} \rangle, \langle \{L\}, \{A\} \rangle, \langle \{L\}, \{S\} \rangle$

2) Candidate pruning:

for all of generated 2-sequences don't have infrequent 1-subsequence, so the pruned result shows below:

$\langle \{U, A\} \rangle, \langle \{U, S\} \rangle, \langle \{U, L\} \rangle, \langle \{A, S\} \rangle, \langle \{A, L\} \rangle, \langle \{S, L\} \rangle,$   
 $\langle \{U\}, \{U\} \rangle, \langle \{U\}, \{A\} \rangle, \langle \{U\}, \{S\} \rangle, \langle \{U\}, \{L\} \rangle, \langle \{S\}, \{S\} \rangle, \langle \{S\}, \{U\} \rangle, \langle \{S\}, \{A\} \rangle, \langle \{S\}, \{L\} \rangle,$   
 $\langle \{A\}, \{A\} \rangle, \langle \{A\}, \{U\} \rangle, \langle \{A\}, \{S\} \rangle, \langle \{A\}, \{L\} \rangle, \langle \{L\}, \{L\} \rangle, \langle \{L\}, \{U\} \rangle, \langle \{L\}, \{A\} \rangle, \langle \{L\}, \{S\} \rangle$

(c) What is the **frequent** 2-sequences based on the results of (b)? (5 Marks)

Candidate elimination:  $\langle \{V, S\} \rangle, \langle \{V, L\} \rangle, \langle \{A, S\} \rangle, \langle \{A, L\} \rangle, \langle \{V\}, \{A\} \rangle, \langle \{V\}, \{S\} \rangle, \langle \{A\}, \{V\} \rangle, \langle \{A\}, \{S\} \rangle, \langle \{S\}, \{V\} \rangle, \langle \{S\}, \{A\} \rangle, \langle \{S\}, \{L\} \rangle, \langle \{L\}, \{V\} \rangle, \langle \{L\}, \{A\} \rangle, \langle \{L\}, \{S\} \rangle, \langle \{V\}, \{U\} \rangle, \langle \{A\}, \{U\} \rangle, \langle \{S\}, \{U\} \rangle, \langle \{A\}, \{L\} \rangle$  are eliminated for their support less than 40%  
⇒ result:  $\langle \{V, A\} \rangle, \langle \{S, L\} \rangle, \langle \{V\}, \{L\} \rangle, \langle \{A\}, \{L\} \rangle$

(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)

1) 3-sequences Candidate Generation by merging:

$\langle \{V, A\}, \{L\} \rangle, \langle \{V\}, \{S, L\} \rangle, \langle \{A\}, \{S, L\} \rangle$

2) Candidate Pruning:

for  $\langle \{V, A\}, \{L\} \rangle$ :  $\langle \{V, A\} \rangle, \langle \{V\}, \{L\} \rangle, \langle \{A\}, \{L\} \rangle$  are all frequent

for  $\langle \{V\}, \{S, L\} \rangle$ :  $\langle \{V\}, \{S\} \rangle$  is infrequent, pruned.

for  $\langle \{A\}, \{S, L\} \rangle$ :  $\langle \{A\}, \{S\} \rangle$  is infrequent, pruned.

(e) What is the **frequent** 3-sequences based on the results of (d)? Please calculate the support. (5 Marks)

for  $\langle \{V, A\}, \{L\} \rangle$ : support = 20% < 40%, eliminated

There were not any Candidate 3-sequences have their support  $\geq$  min-support.