

**BACS3013 Data Science**

**Tutorial 6 (Descriptive Analytics - part 1)**

Q1. Find all association rules in the following database:  
(minsup = 2, minconf = 70% )

TID	Items
1	a, b, c
2	b, c, d, e
4	c, d
3	a, b, d
5	a, b, c

**Ans:**

**Firstly, find all itemsets with support  $\geq 2$ ;**

$C_1$		$C_2$		$C_3$	
set	supp	set	supp	set	supp
{a}	3	{a,b}	3	{a,b,c}	2
{b}	4	{a,c}	2	<del>{b,c,d}</del>	<del>1</del>
{c}	4	<del>{a,d}</del>	<del>1</del>		
{d}	3	{b,c}	3		
<del>{e}</del>	<del>1</del>	{b,d}	2		
		{c,d}	2		

$F = \{\{a\}, \{b\}, \{c\}, \{d\}, \{a,b\}, \{a,c\}, \{b,c\}, \{b,d\}, \{c,d\}, \{a,b,c\}\}$

**Then, split all frequent itemsets in all possible ways**

BACS3013 Data Science

set	supp
{}	5
{a}	3
{b}	4
{c}	4
{d}	3
{a,b}	3
{a,c}	2
{b,c}	3
{b,d}	2
{c,d}	2
{a,b,c}	2

**(Trivial rules:  $X \Rightarrow \{\}$  100%)**

<b><math>ac \Rightarrow b</math></b>	<b>100%</b>
<b><math>a \Rightarrow b</math></b>	<b>100%</b>
<b><math>b \Rightarrow a</math></b>	<b>75%</b>
<b><math>b \Rightarrow c</math></b>	<b>75%</b>
<b><math>c \Rightarrow b</math></b>	<b>75%</b>
<b><math>\{\} \Rightarrow b</math></b>	<b>80%</b>
<b><math>\{\} \Rightarrow c</math></b>	<b>80%</b>

BACS3013 Data Science

Q2. For the confidence of association rules, a monotonicity principle can be stated between rules that are based on the same itemset ( $X \Rightarrow Y$  is based on  $X \cup Y$ ). Explain how can it be exploited?

Ans:

$$\begin{aligned} \text{Let } I = X \cup Y = X' \cup Y', \text{ and let } X' \subseteq X \\ \text{Then:} \quad \sup(X \Rightarrow Y) = \sup(X' \Rightarrow Y') \\ \text{conf}(X \Rightarrow Y) \geq \text{conf}(X' \Rightarrow Y') \end{aligned}$$

**Example:**

$$\begin{aligned} \text{conf}(\{a,b\} \Rightarrow \{c\}) &\geq \text{conf}(\{a\} \Rightarrow \{b,c\}) \\ &\geq \text{conf}(\{\} \Rightarrow \{a,b,c\}) \end{aligned}$$

$$\begin{aligned} \sup(X \Rightarrow Y) &= \sup(X \cup Y) = \sup(I) \\ &= \sup(X' \cup Y') = \sup(X' \Rightarrow Y') \end{aligned}$$

$$\begin{aligned} \text{conf}(X \Rightarrow Y) &= \sup(X \cup Y) / \sup(X) \\ &\geq \sup(X \cup Y) / \sup(X') \quad (X' \subseteq X) \\ &= \sup(X' \cup Y') / \sup(X) \quad (X \cup Y = X' \cup Y') \\ &= \text{conf}(X' \Rightarrow Y') \end{aligned}$$

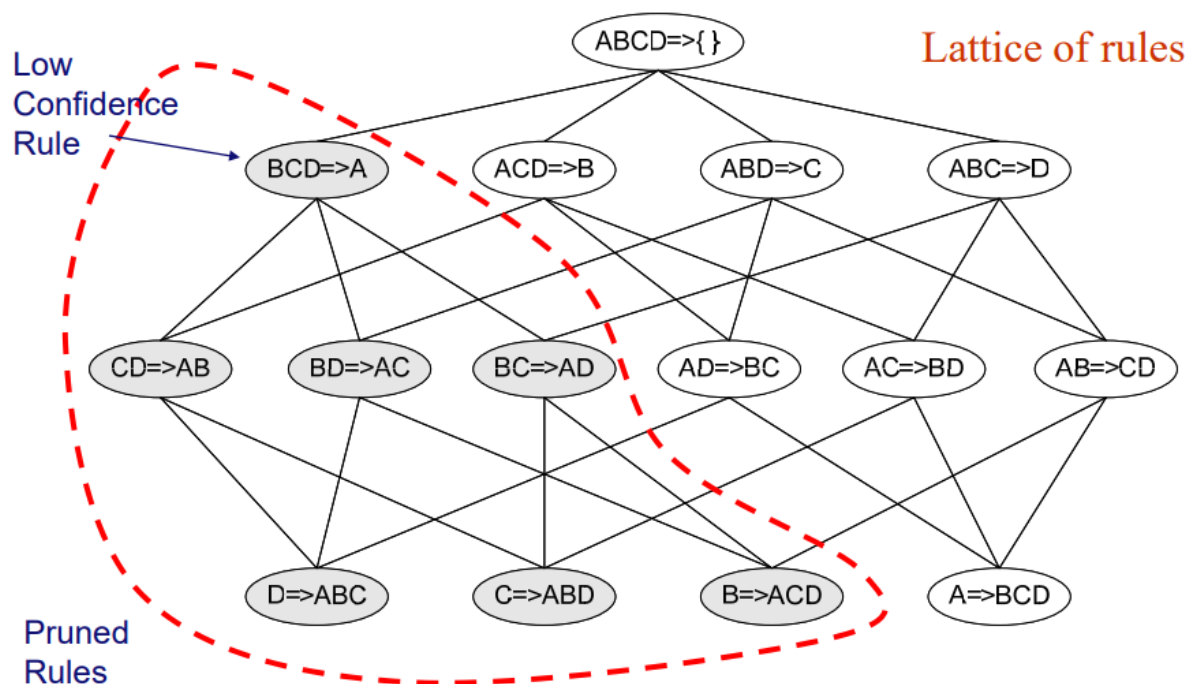
**Note:**

**When splitting  $I$  into  $I = X \cup Y$ , first consider singleton  $Y$**

**Only consider  $X \Rightarrow Y$  if:**

**for all  $y \in Y$ ,  $X \cup \{y\} \Rightarrow Y \setminus \{y\}$  was confident**

**BACS3013 Data Science**



Q3. A transaction database is shown below:

TID	Products
1	A, B, E
2	B, D
3	B, C
4	A, B, D
5	A, C
6	B, C
7	A, C
8	A, B, C, E
9	A, B, C

**BACS3013 Data Science**

- a. Why is every subset of a frequent set frequent?
- b. Given frequent set  $\{A, B, E\}$ , what are the possible association rules?
- c. Given frequent set  $\{A, B, E\}$ , what association rules have  $\text{minsup}=2$  and  $\text{minconf}=50\%$ ?

**Ans:**

- a. Suppose  $\{A, B\}$  is frequent. Since each occurrence of  $A, B$  includes both  $A$  and  $B$ , then both  $A$  and  $B$  must also be frequent.

b.

- $A \Rightarrow B, E$
- $A, B \Rightarrow E$
- $A, E \Rightarrow B$
- $B \Rightarrow A, E$
- $B, E \Rightarrow A$
- $E \Rightarrow A, B$
- $\_ \Rightarrow A, B, E$  (empty rule), or  $\text{true} \Rightarrow A, B, E$

c.

BACS3013 Data Science

$A, B \Rightarrow E : \text{conf} = 2/4 = 50\%$

$A, E \Rightarrow B : \text{conf} = 2/2 = 100\%$

$B, E \Rightarrow A : \text{conf} = 2/2 = 100\%$

$E \Rightarrow A, B : \text{conf} = 2/2 = 100\%$

Don't qualify

$A \Rightarrow B, E : \text{conf} = 2/6 = 33\% < 50\%$

$B \Rightarrow A, E : \text{conf} = 2/7 = 28\% < 50\%$

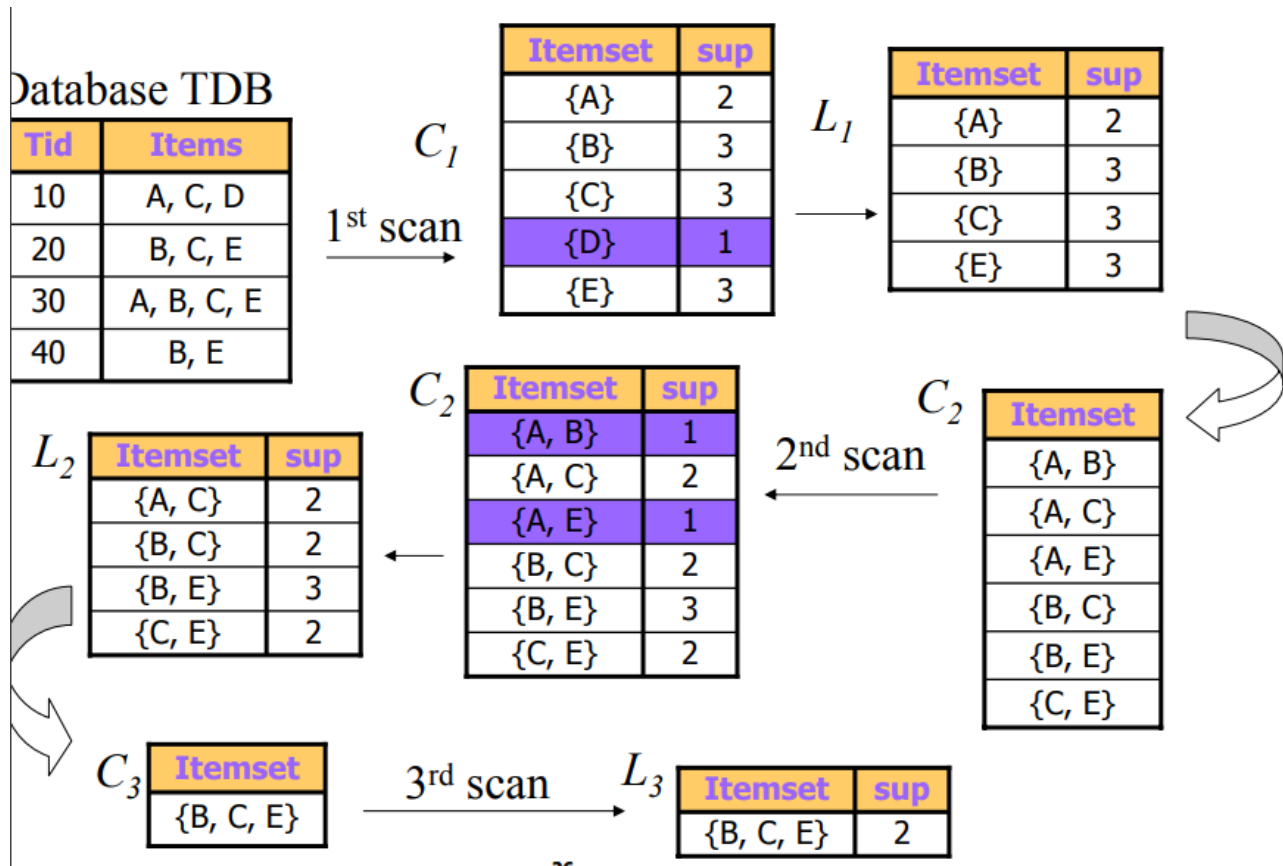
$\_ \Rightarrow A, B, E : \text{conf} = 2/9 = 22\% < 50\%$

Q4. By using Apriori algorithm, mine the association rules for the following transaction:

Transaction-id	Items bought
10	A, B, C
20	A, C
30	A, D
40	B, E, F

Ans:

**BACS3013 Data Science**



Q5. The following is an example of customer purchase transaction data set.

CID	TID	Date	Items Purchased
1	1	01/01/2001	10,20
1	2	01/02/2001	10,30,50,70
1	3	01/03/2001	10,20,30,40
2	4	01/03/2001	20,30
2	5	01/04/2001	20,40,70
3	6	01/04/2001	10,30,60,70

**BACS3013 Data Science**

3	7	01/05/2001	10,50,70
4	8	01/05/2001	10,20,30
4	9	01/06/2001	20,40,60
5	10	01/11/2001	10,20,30,60

Note: CID = Customer ID and TID = Transactions ID

- a) Calculate the support, confidence and lift of the following association rule. Indicate if the items in the association rule are independent of each other or have negative or positive impacts on each other.  
 $\{10\} \rightarrow \{50,70\}$
- b) The following is the list of large two item sets. Show the steps to apply the Apriori property to generate and prune the candidates for large three itemsets. Describe how the Apriori property is used in the steps. Give the final list of candidate large three item sets.  
 $\{10,20\} \{10,30\} \{20,30\} \{20,40\}$
- c) Does customer 1 support the sequence  $\langle \{20\} \{50,70\} \{10\} \rangle$ ? Justify your answer.
- d) Calculate the support of  $\langle \{10\}, \{30\} \rangle$ .
- e) Based on the types of association rules discussed in class, identify which type(s) of rules  $\{10\} \rightarrow \{50,70\}$  is.

**Ans:**

**a)**

**Support** =  $\text{Support}(\{10,50,70\}) = 2/10 = 20\%$

**Confidence** =  $\text{Support}(\{10,50,70\}) / \text{Support}(\{10\}) = 0.2/0.7 = 2/7 = 29\%$

**Lift** =  $\text{Confidence} / \text{Support}(\{50,70\}) = 2/7 / 0.2 = 10/7 = 1.43 > 1$

Since lift is larger than 1, it's a positive rule.

**b)**

$\{10,20\} \{10,30\} \{20,30\} \{20,40\}$

**\*\*\*O:** describe how the apriori property is used to decide which 2 large item sets are joined together and to determine which 3 item set should be pruned.

**Join:**  $\{10,20,30\} \{20,30,40\}$

**Prune:**  $\{10,20,30\}$  ( $\{20,30,40\}$  is pruned)

**Final list:**  $\{10,20,30\}$

**d)**



**BACS3013 Data Science**

The sequence of customer 1 is:

$\langle \{10,20\} \{10,30,50,70\} \{10,20,30,40\} \rangle$

Since  $\{20\} \subseteq \{10,20\}$ ,  $\{50,70\} \subseteq \{10,30,50,70\}$ , and  $\{10\} \subseteq \{10,20,30,40\}$ ,

$\langle \{20\} \{50,70\} \{10\} \rangle$  is contained in the sequence of customer 1. Therefore, customer 1 supports sequence  $\langle \{20\} \{50,70\} \{10\} \rangle$ .

e)

Only customer 1 supports the sequence  $\langle \{10\} \{30\} \rangle$  and there are 5 customers, therefore,

Support =  $1/5 = 20\%$

f)

The association rule  $\{10\} \rightarrow \{50,70\}$  is a single-level, single-dimensional and Boolean association rule.