

COMP1942 Exploring and Visualizing Data (Spring Semester 2019)

Midterm Examination (Answer Sheet)

Date: 8 April, 2019 (Monday)

Time: 10:40am-11:50am

Duration: 1 hour 10 minutes

Student ID: \_\_\_\_\_

Student Name: \_\_\_\_\_

Seat No. : \_\_\_\_\_

Instructions:

- (1) Please answer **all** questions in this paper.
- (2) The total marks are 100.
- (3) You can use a calculator.

# Answer Sheet

Question	Full Mark	Mark
Q1	20	
Q2	20	
Q3	20	
Q4	20	
Q5	20	
Total	100	

**Q1 (20 Marks)**

(a) (i)

No.

Consider the following example.

<b>B</b>	<b>C</b>
1	1
1	1
1	1

In Step 1,  $\{B, C\}$  and  $\{B\}$  are in  $S_1$   
 since  $\text{supp}(\{B, C\}) \geq 3$  and  $\text{supp}(\{B\}) \geq 3$

In Step 2,  $B \rightarrow C$  is generated since  $\text{supp}(\{B, C\}) / \text{supp}(\{B\}) = 100\% \geq 50\%$

Thus,  $B \rightarrow C$  is in  $S_2$

Note that

$$\begin{aligned} \text{supp}(B \rightarrow C) &= \text{supp}(\{B, C\}) \\ &= 3 \\ &< 4 \end{aligned}$$

In conclusion,  $B \rightarrow C$  is in  $S_2$  but  $\text{supp}(B \rightarrow C) < 4$

**Q1 (continued)**

(a) (ii)

Yes.

Since  $B \rightarrow C$  is in  $S_0$ ,

$$\text{conf}(B \rightarrow C) \geq 50\%$$

$$\text{supp}(\{B, C\}) / \text{supp}(\{B\}) \geq 50\%$$

Since  $B \rightarrow C$  is in  $S_0$ ,

$$\text{supp}(B \rightarrow C) \geq 4$$

Since  $\text{supp}(\{B, C\}) = \text{supp}(B \rightarrow C)$ ,

$$\text{supp}(\{B, C\}) \geq 4$$

Thus,  $\{B, C\}$  is in  $S_1$ .Since  $\text{supp}(\{B, C\}) \geq 4$ ,

$$\text{supp}(\{B\}) \geq 4$$

Thus,  $\{B\}$  is in  $S_1$ Since  $\{B\}$  is in  $S_1$ ,and  $\{B, C\}$  is in  $S_1$ ,

Step 2 must consider

 $\{B\}$  and  $\{B, C\}$  together, andgenerate  $B \rightarrow C$  (since  $\text{supp}(\{B, C\}) / \text{supp}(\{B\}) \geq 50\%$ ) $B \rightarrow C$  is in  $S_2$ .

**Q1 (continued)**

(b)(i)

Yes.

Since " $B \rightarrow C$ " is in  $S_2$ ,

we know that

we have to calculate  $\text{supp}(\{B, C\})/\text{supp}(\{B\})$  in Step (\*)

In other words,

$\{B, C\}$  and  $\{B\}$  are in  $S_1$

which means that

$$\text{supp}(\{B, C\}) \geq 4 \text{ and}$$

$$\text{supp}(\{B\}) \geq 4$$

Since  $\text{supp}(B \rightarrow C) = \text{supp}(\{B, C\})$ ,

$$\text{supp}(B \rightarrow C) \geq 4$$

Thus,

$$\text{supp}(B \rightarrow C) \geq 4$$

**Q1 (continued)**

(b) (ii)

No.

Consider the following example.

<b>B</b>	<b>C</b>
1	1
1	1
1	1
1	1
1	0
1	0
1	0

 $B \rightarrow C$  is in  $S_0$ (since  $\text{conf}(B \rightarrow C) = 57.14\%$  and  $\text{supp}(B \rightarrow C) = 4$ )Since  $\text{supp}(\{B, C\}) = 4$  and  $\text{supp}(\{B\}) = 7$ ,both  $\{B, C\}$  and  $\{B\}$  are in  $S_1$ .However, in Step 2,  $B \rightarrow C$  is not generated in the output set  $S_2$  because $\text{supp}(\{B, C\}) / \text{supp}(\{B\}) = 57.14\%$  which is smaller than 60%.

**Q2 (20 Marks)**

(a)

$$L_1 = \{\{A\}, \{C\}, \{D\}, \{E\}\}$$

Large 2-itemset Generation:

Join Step/Prune Step

$$C_2 = \{\{A, C\}, \{A, D\}, \{A, E\}, \{C, D\}, \{C, E\}, \{D, E\}\}$$

Counting Step

$$L_2 = \{\{A, C\}, \{A, D\}, \{A, E\}, \{C, E\}, \{D, E\}\}$$

Large 3-itemset Generation:

Join Step

$$C_3 = \{\{A, C, D\}, \{A, C, E\}, \{A, D, E\}\}$$

Prune Step

$$C_3 = \{\{A, C, E\}, \{A, D, E\}\}$$

Counting Step

$$L_3 = \{\{A, C, E\}, \{A, D, E\}\}$$

Large 4-item set Generation:

Join Step

$$C_4 = \{\}$$

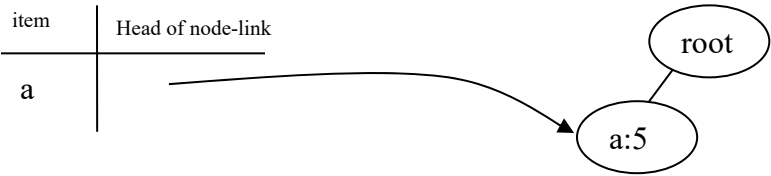
$$\text{Large itemsets} = L_1 \cup L_2 \cup L_3$$

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

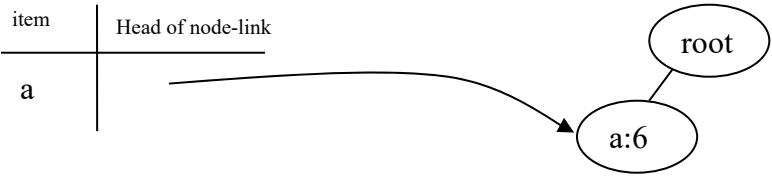
**Q2 (Continued)**

Q2 (Continued)

(b)(i)



(ii)



(iii)



(iv)

{a}, {b}, {c}, {a, c}, {a, b}



**Q3 [20 Marks]**

(a)(i)

$$\text{Info}(T) = 1 - 0.5^2 - 0.5^2 = 0.5$$

For attribute Age,

$$\text{Info}(T_{\text{young}}) = 1 - 0.5^2 - 0.5^2 = 0.5$$

$$\text{Info}(T_{\text{old}}) = 1 - 0.5^2 - 0.5^2 = 0.5$$

$$\text{Info}(\text{Age}, T) = \frac{1}{2} \text{Info}(T_{\text{young}}) + \frac{1}{2} \text{Info}(T_{\text{old}}) = 0.5$$

$$\text{Gain}(\text{Age}, T) = \text{Info}(T) - \text{Info}(\text{Age}, T) = 0$$

For attribute Gender,

$$\text{Info}(T_{\text{male}}) = 1 - 1^2 - 0^2 = 0$$

$$\text{Info}(T_{\text{female}}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.4444$$

$$\text{Info}(\text{Gender}, T) = \frac{1}{4} \text{Info}(T_{\text{male}}) + \frac{3}{4} \text{Info}(T_{\text{female}}) = 0.3333$$

$$\text{Gain}(\text{Gender}, T) = \text{Info}(T) - \text{Info}(\text{Gender}, T) = 0.1667$$

For attribute MMR\_Vaccine,

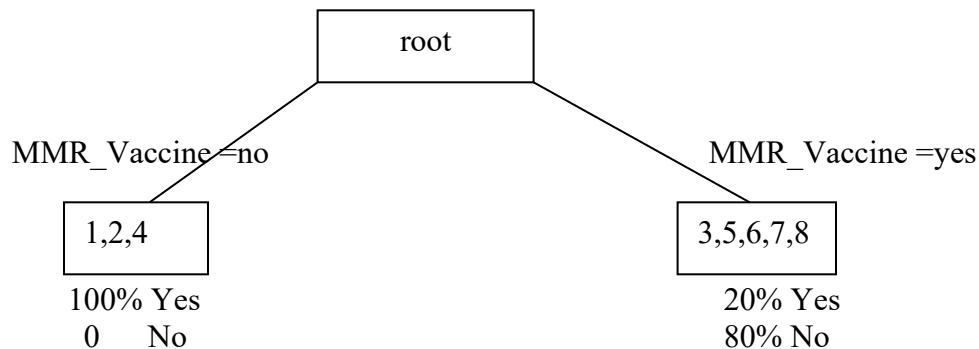
$$\text{Info}(T_{\text{no}}) = 1 - 1^2 - 0^2 = 0$$

$$\text{Info}(T_{\text{yes}}) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 = 0.32$$

$$\text{Info}(\text{MMR\_Vaccine}, T) = \frac{3}{8} \text{Info}(T_{\text{no}}) + \frac{5}{8} \text{Info}(T_{\text{yes}}) = 0.2$$

$$\text{Gain}(\text{MMR\_Vaccine}, T) = \text{Info}(T) - \text{Info}(\text{MMR\_Vaccine}, T) = 0.3$$

We choose attribute MMR\_Vaccine for Splitting:



**Q3 (Continued)**

**Q3 (Continued)**

(a) (ii)

It is likely that he will not have measles.

(b)

Differences:

The definition of the gain used in C4.5 is different from that used in ID3.

The gain used in C4.5 is equal to the gain used in ID3 divided by SplitInfo.

The reason why there is a difference is described as follows.

In ID3, there is a higher tendency to choose an attribute containing more values (e.g., attribute identifier and attribute HKID). Thus, splitInfo in C4.5 is used to penalize an attribute containing more values. If this value is larger, the penalty is larger.

**Q4 [20 Marks]**

(a)

Yes.

Cluster 1:  $\{x_1, x_2, x_4, x_5, x_6\}$

Cluster 2:  $\{x_3, x_7, x_8, x_9, x_{10}\}$

(b) (i)

Yes.

Cluster 1:  $\{x_2, x_5, x_6\}$

Cluster 2:  $\{x_3, x_7, x_8, x_9, x_{10}\}$

Cluster 3:  $\{x_4\}$

Cluster 4:  $\{x_1\}$

(ii)

Yes.

Cluster 1:  $\{x_1, x_2, x_4, x_5, x_6\}$

Cluster 2:  $\{x_3, x_7, x_8, x_9, x_{10}\}$

**Q4 (Continued)**

(b) (iii)

No.

This is because in the dendrogram, we have to specify the distance between 2 clusters.

However, when we use the centroid linkage as a distance measurement between 2 clusters, we have to know the coordinate of each point (and thus the mean/center of all points in each cluster), which could not be found in the information given.

**Q5 [20 Marks]**

(a)(i)

- Make initial guesses for the means  $m_1, m_2, \dots, m_k$
- Until Interrupted
  - Acquire the next example  $x$
  - If  $m_i$  is closest to  $x$ ,
    - replace  $m_i$  by  $m_i + a(x - m_i)$

(ii)

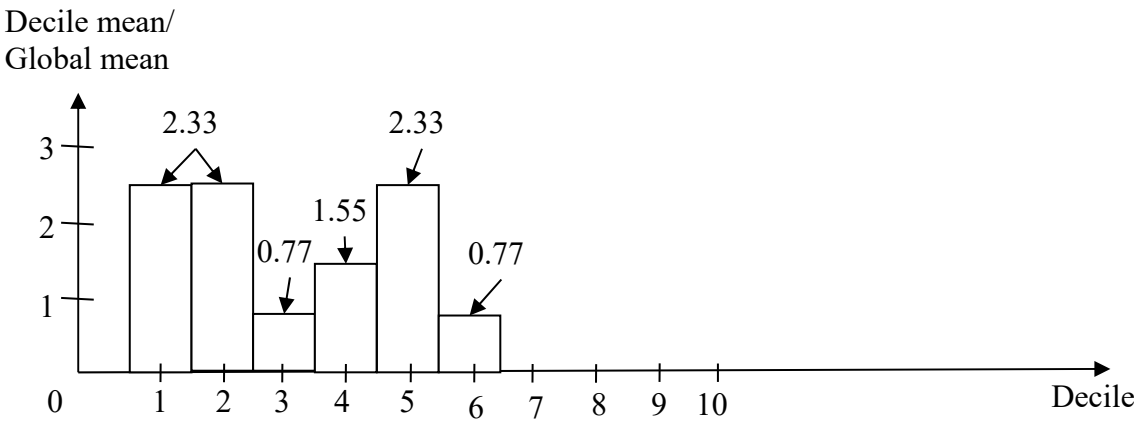
$$\begin{aligned}
 m_n &= m_{n-1} + a(x_n - m_{n-1}) \\
 &= (1-a)m_{n-1} + ax_n \\
 &= (1-a)[(1-a)m_{n-2} + ax_{n-1}] + ax_n \\
 &= (1-a)^2 m_{n-2} + (1-a)ax_{n-1} + ax_n \\
 &= (1-a)^2 [(1-a)m_{n-3} + ax_{n-2}] + (1-a)ax_{n-1} + ax_n \\
 &= (1-a)^3 m_{n-3} + (1-a)^2 ax_{n-2} + (1-a)ax_{n-1} + ax_n \\
 &= \dots \\
 &= (1-a)^n m_0 + \sum_{p=1}^n (1-a)^{n-p} ax_p
 \end{aligned}$$

$$X = (1-a)^n$$

$$Y = (1-a)^{n-p} a$$

**Q5 (Continued)**  
(b)

Yes. The chart is shown as follows.



**End of Answer Sheet**