

COMP1942 Exploring and Visualizing Data (Spring Semester 2024)
Homework 1 Solution
Full Mark: 100 Marks

Q1 [20 Marks]

(a)

Item	Count
AngusBeefBurger	4
CrabRoeNoodle	4
LimeJuice	3
MulberryJuice	1
RoastDuck	4
ScallopCongee	4

$L_1 = \{\{\text{AngusBeefBurger}\}, \{\text{CrabRoeNoodle}\}, \{\text{LimeJuice}\}, \{\text{RoastDuck}\}, \{\text{ScallopCongee}\}\}$

Generate L_2 :

Join step/prune step:

$C_2 = \{\{\text{AngusBeefBurger}, \text{CrabRoeNoodle}\}, \{\text{AngusBeefBurger}, \text{LimeJuice}\},$
 $\{\text{AngusBeefBurger}, \text{RoastDuck}\}, \{\text{AngusBeefBurger}, \text{ScallopCongee}\},$
 $\{\text{CrabRoeNoodle}, \text{LimeJuice}\}, \{\text{CrabRoeNoodle}, \text{RoastDuck}\}, \{\text{CrabRoeNoodle}, \text{ScallopCongee}\},$
 $\{\text{LimeJuice}, \text{RoastDuck}\}, \{\text{LimeJuice}, \text{ScallopCongee}\},$
 $\{\text{RoastDuck}, \text{ScallopCongee}\}\}$

Counting step:

$L_2 = \{\{\text{AngusBeefBurger}, \text{CrabRoeNoodle}\}, \{\text{AngusBeefBurger}, \text{LimeJuice}\}, \{\text{AngusBeefBurger},$
 $\text{ScallopCongee}\},$
 $\{\text{CrabRoeNoodle}, \text{RoastDuck}\}, \{\text{CrabRoeNoodle}, \text{ScallopCongee}\},$
 $\{\text{LimeJuice}, \text{RoastDuck}\},$
 $\{\text{RoastDuck}, \text{ScallopCongee}\}$
 $\}$

Generate L_3 :

Join step:

$C_3 = \{\{\text{AngusBeefBurger}, \text{CrabRoeNoodle}, \text{LimeJuice}\}, \{\text{AngusBeefBurger}, \text{CrabRoeNoodle},$
 $\text{ScallopCongee}\},$
 $\{\text{AngusBeefBurger}, \text{LimeJuice}, \text{ScallopCongee}\},$
 $\{\text{CrabRoeNoodle}, \text{RoastDuck}, \text{ScallopCongee}\}\}$

Prune step:

$C_3 = \{\{\text{AngusBeefBurger}, \text{CrabRoeNoodle}, \text{ScallopCongee}\}, \{\text{CrabRoeNoodle}, \text{RoastDuck},$
 $\text{ScallopCongee}\}\}$

Counting step:

$L_3 = \{\{\text{CrabRoeNoodle}, \text{RoastDuck}, \text{ScallopCongee}\}\}$

Large itemsets = $L_1 \cup L_2 \cup L_3$

$= \{\{\text{AngusBeefBurger}\}, \{\text{CrabRoeNoodle}\}, \{\text{LimeJuice}\}, \{\text{RoastDuck}\}, \{\text{ScallopCongee}\},$
 $\{\text{AngusBeefBurger}, \text{CrabRoeNoodle}\}, \{\text{AngusBeefBurger}, \text{LimeJuice}\},$
 $\{\text{AngusBeefBurger}, \text{ScallopCongee}\},$
 $\{\text{CrabRoeNoodle}, \text{RoastDuck}\}, \{\text{CrabRoeNoodle}, \text{ScallopCongee}\},$
 $\{\text{LimeJuice}, \text{RoastDuck}\}, \{\text{RoastDuck}, \text{ScallopCongee}\}\}$

{LimeJuice, RoastDuck},
 {RoastDuck, ScallopCongee},
 {CrabRoeNoodle, RoastDuck, ScallopCongee } }

(b) Association rules

{RoastDuck, ScallopCongee} → CrabRoeNoodle
 {CrabRoeNoodle, RoastDuck} → ScallopCongee

Q2 [20 Marks]

(a)

Item	Freq
a	3
b	3
c	5
d	5
e	1
f	6
g	1
h	1
i	1
j	1
k	1
l	1
m	1
n	1
o	1
p	1
q	1

Freq items:

Item	Freq
a	3
b	3
c	5
d	5
f	6

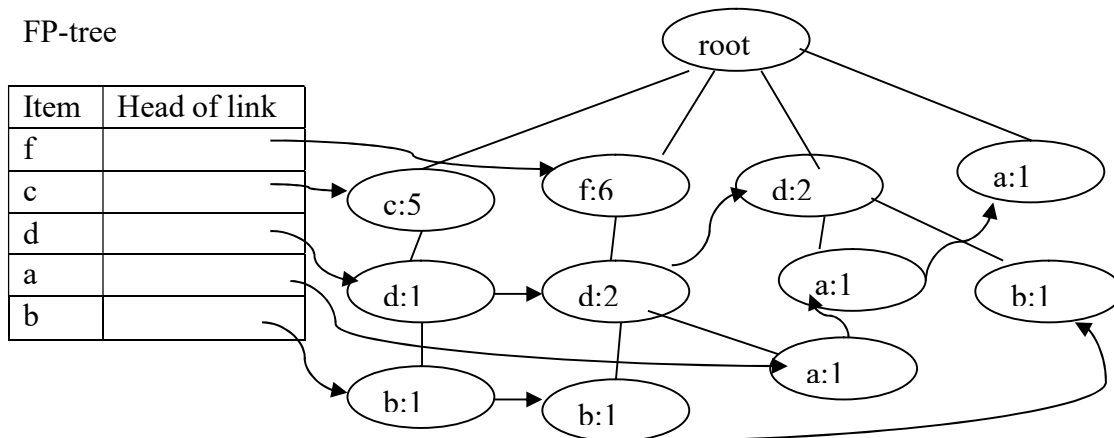
Sorted Freq items:

Item	Freq
f	6
c	5
d	5
a	3
b	3

Ordered freq items

TID	Items bought	(ordered) freq items
1	b,c,d,p	c,d,b
2	f,j,q	f
3	c,i	c
4	a,d	d,a
5	c,m	c
6	b,d,f	f,d,b
7	a,d,f	f,d,a
8	a,l	a
9	c,g	c
10	c,k	c
11	f,n,o	F
12	e,f	f
13	f,h	f
14	b,d	d,b

FP-tree



Conditional FP-tree on “b”

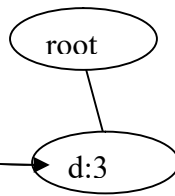
count (b)=3

$$\begin{array}{ll} (c:1,d:1,b:1) & (b:1,d:1) \\ (f:1,d:1,b:1) & \Rightarrow (b:1,d:1) \\ (d:1,b:1) & (b:1,d:1) \end{array}$$

Item	Freq
f	1
c	1
d	3
a	0
b	3

Item	Freq
b	3
d	3

Item	Head
d	



{b,d}:3

Conditional FP-tree on “a”

(f:1,d:1,a:1) (a:1,d:1)
 (d:1,a:1) ⇒ (a:1,d:1)
 (a:1) (a:1)

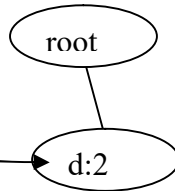
count (a)=3

Item	Freq
f	1
c	0
d	2
a	3
b	0

↓

Item	Freq
a	3
d	2

Item	Head
d	



{a,d}:2

Conditional FP-tree on “d”

(c:1,d:1) (d:1)
 (f:2,d:2) ⇒ (d:2,f:2)
 (d:2) (d:2)

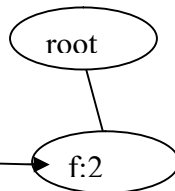
count (d)=5

Item	Freq
f	2
c	1
d	5
a	0
b	0

↓

Item	Freq
d	5
f	2

Item	Head
f	



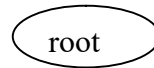
{d,f}:2

Conditional FP-tree on “c” $(c:5) \Rightarrow (c:5)$

Item	freq
c	5



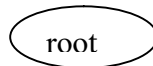
Item	freq
c	5

count (c)=5**Conditional FP-tree on “f”** $(f:6) \Rightarrow (f:6)$

Item	freq
f	6



Item	freq
f	6

count (f)=6**Freq itemsets**

= { {b}, {b,d},
 {a}, {a,d},
 {d}, {d,f},
 {c},
 {f} }

(b) association rules

$\{b\} \rightarrow \{d\}$
 $\{d\} \rightarrow \{b\}$
 $\{a\} \rightarrow \{d\}$

Q3 [20 Marks]

(a)

(i)

cluster 1:

mean=(16.5,18.25)

points={ x_1, x_3, x_5, x_8 }

cluster 2:

mean=(6.25,10)

points={ x_2, x_4, x_6, x_7 }

(ii)

cluster1:

mean=(4.33,12.33)

points={ x_2, x_4, x_6 }

cluster2:

mean=(16.5,18.25)

points={ x_1, x_3, x_5, x_8 }

cluster 3:

mean=(12,3)

points={ x_7 }

(b)

Advantages:

1. It is easy to implement.
2. It is easy to understand.

Disadvantages:

1. It is difficult to determine the value of k because we do not know the no. of the clusters.
2. K-means is sensitive to the initial guess of the means.

Q4 [20 Marks]

	1	2	3	4	5	6	7	8
1	0							
2	11	0						
3	5	13	0					
4	12	2	14	0				
5	7	7	1	18	0			
6	13	4	15	5	20	0		
7	9	15	12	16	15	19	0	
8	11	20	12	21	17	22	30	0

$$D(1,*)=13$$

$$D(2,*)=20$$

$$D(3,*)=15$$

$$D(4,*)=21$$

$$D(5,*)=20$$

$$D(6,*)=22$$

$$D(7,*)=30$$

$$D(8,*)=30$$

$$A=\{8\}$$

$$B=\{1,2,3,4,5,6,7\}$$

$$D(1,A)=11 \quad D(1,B)=13 \quad \Delta_1=2$$

$$D(2,A)=20 \quad D(2,B)=17 \quad \Delta_2=-3$$

$$D(3,A)=12 \quad D(3,B)=15 \quad \Delta_3=3$$

$$D(4,A)=21 \quad D(4,B)=18 \quad \Delta_4=-3$$

$$D(5,A)=17 \quad D(5,B)=20 \quad \Delta_5=3$$

$$D(6,A)=22 \quad D(6,B)=20 \quad \Delta_6=-2$$

$$D(7,A)=30 \quad D(7,B)=19 \quad \Delta_7=-11$$

$$A=\{3,8\}$$

$$B=\{1,2,4,5,6,7\}$$

$$D(1,A)=11 \quad D(1,B)=13 \quad \Delta_1=2$$

$$D(2,A)=20 \quad D(2,B)=17 \quad \Delta_2=-3$$

$$D(4,A)=21 \quad D(4,B)=18 \quad \Delta_4=-3$$

$$D(5,A)=17 \quad D(5,B)=20 \quad \Delta_5=3$$

$$D(6,A)=22 \quad D(6,B)=20 \quad \Delta_6=-2$$

$$D(7,A)=30 \quad D(7,B)=19 \quad \Delta_7=-11$$

$$A=\{3,5,8\}$$

$$B=\{1,2,4,6,7\}$$

$$D(1,A)=11 \quad D(1,B)=13 \quad \Delta_1=2$$

$$D(2,A)=20 \quad D(2,B)=15 \quad \Delta_2=-5$$

$D(4,A)=21$	$D(4,B)=16$	$\Delta_4 = -5$
$D(6,A)=22$	$D(6,B)=19$	$\Delta_6 = -3$
$D(7,A)=30$	$D(7,B)=19$	$\Delta_7 = -11$

$A=\{1,3,5,8\}$
 $B=\{2,4,6,7\}$

$D(2,A)=20$	$D(2,B)=15$	$\Delta_2 = -5$
$D(4,A)=21$	$D(4,B)=16$	$\Delta_4 = -5$
$D(6,A)=22$	$D(6,B)=19$	$\Delta_6 = -3$
$D(7,A)=30$	$D(7,B)=19$	$\Delta_7 = -11$

Stop!
 $A=\{1,3,5,8\}$
 $B=\{2,4,6,7\}$

So, there are two clusters:

Cluster 1: data points 1, 3, 5, 8

Cluster 2: data points 2, 4, 6, 7

The distance between 2 clusters is 30.

Q5 [20 Marks]

(a)(i)

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

For attribute HasMacBook,

$$\text{Info}(T_{no}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{yes}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(\text{HasMacBook}, T) = \frac{1}{2} \text{Info}(T_{no}) + \frac{1}{2} \text{Info}(T_{yes}) = 1$$

$$\text{SplitInfo}(\text{HasMacBook}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Gain}(\text{HasMacBook}, T) = \frac{1-1}{1} = 0$$

For attribute Income,

$$\text{Info}(T_{high}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{medium}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{low}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(\text{Income}, T) = \frac{1}{2} \text{Info}(T_{high}) + \frac{1}{4} \text{Info}(T_{medium}) + \frac{1}{4} \text{Info}(T_{low}) = 1$$

$$\text{SplitInfo}(\text{Income}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1.5$$

$$\text{Gain}(\text{Income}, T) = \frac{1-1}{1.5} = 0$$

For attribute Age,

$$\text{Info}(T_{middle}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{old}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

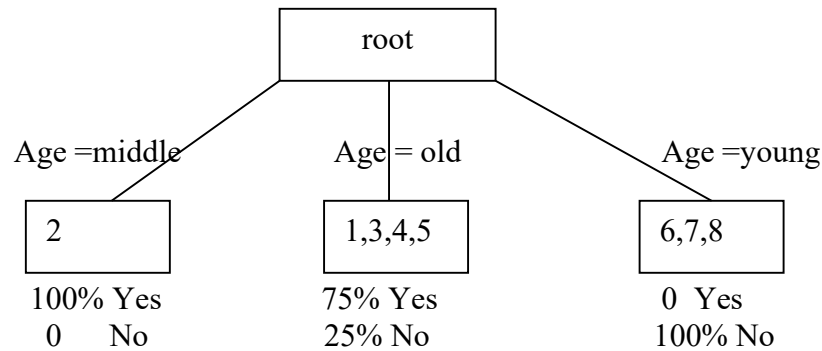
$$\text{Info}(T_{young}) = -0 \log 0 - 1 \log 1 = 0$$

$$\text{Info}(\text{Age}) = \frac{1}{8} \text{Info}(T_{middle}) + \frac{1}{2} \text{Info}(T_{old}) + \frac{3}{8} \text{Info}(T_{young}) = 0.405$$

$$\text{SplitInfo}(\text{Age}) = -\frac{1}{8} \log \frac{1}{8} - \frac{1}{2} \log \frac{1}{2} - \frac{3}{8} \log \frac{3}{8} = 1.4056$$

$$\text{Gain}(\text{Age}, T) = \frac{1-0.405}{1.4056} = 0.4233$$

We choose attribute Age for Splitting:



Consider the node for “Age=old”

$$\text{Info}(T) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

For attribute HasMacBook,

$$\text{Info}(T_{no}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{yes}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(\text{HasMacBook}, T) = \frac{1}{2} \text{Info}(T_{no}) + \frac{1}{2} \text{Info}(T_{yes}) = 0.5$$

$$\text{SplitInfo}(\text{HasMacBook}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Gain}(\text{HasMacBook}, T) = \frac{0.8113 - 0.5}{1} = 0.3113$$

For attribute Income,

$$\text{Info}(T_{\text{high}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{\text{medium}}) = -1 \log 1 - 0 \log 0 = 0$$

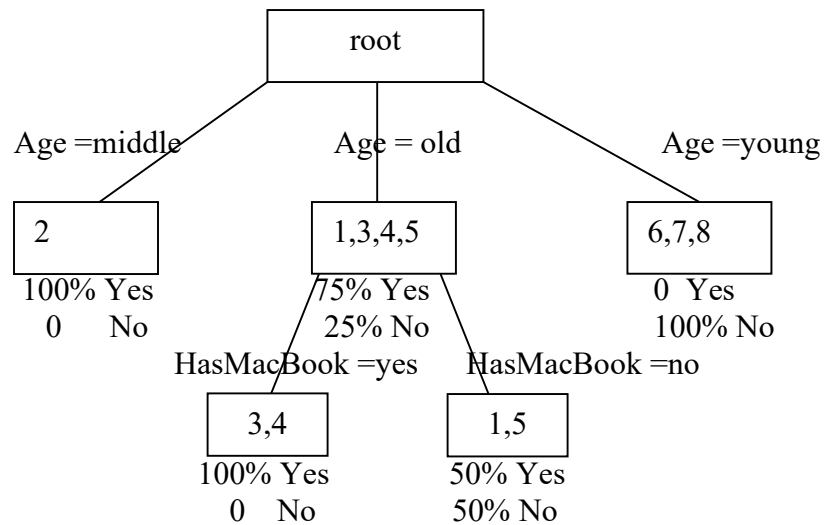
$$\text{Info}(T_{\text{low}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(\text{Income}) = \frac{1}{2} \text{Info}(T_{\text{high}}) + \frac{1}{4} \text{Info}(T_{\text{medium}}) + \frac{1}{4} \text{Info}(T_{\text{low}}) = 0.5$$

$$\text{SplitInfo}(\text{Income}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1.5$$

$$\text{Gain}(\text{Income}, T) = \frac{0.8113 - 0.5}{1.5} = 0.2075$$

We choose attribute HasMacBook for Splitting:



(ii) It is very likely that this customer will buy “Apple Vision Pro”.

(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.