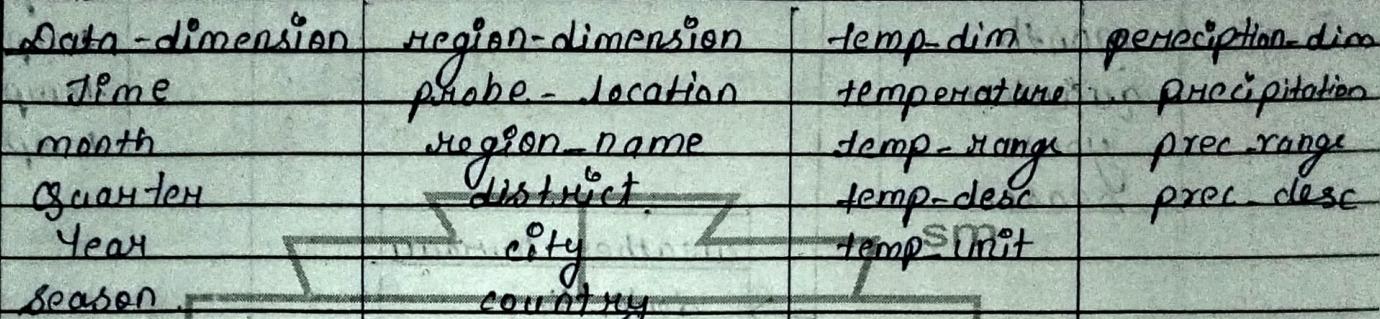


## ASSIGNMENT NO. 1

Q1) Data warehouse for weather bureau  
 → I.P.D:



SHAH & ANCHOR

Fact: region-map(), count(), area(), Pressure()

## STAR SCHEMA

region-dim  
 region location  
 region name  
 district  
 country  
 city  
 state

precipitation dim  
 prec-key  
 precipitation  
 prec-range  
 prec-desc

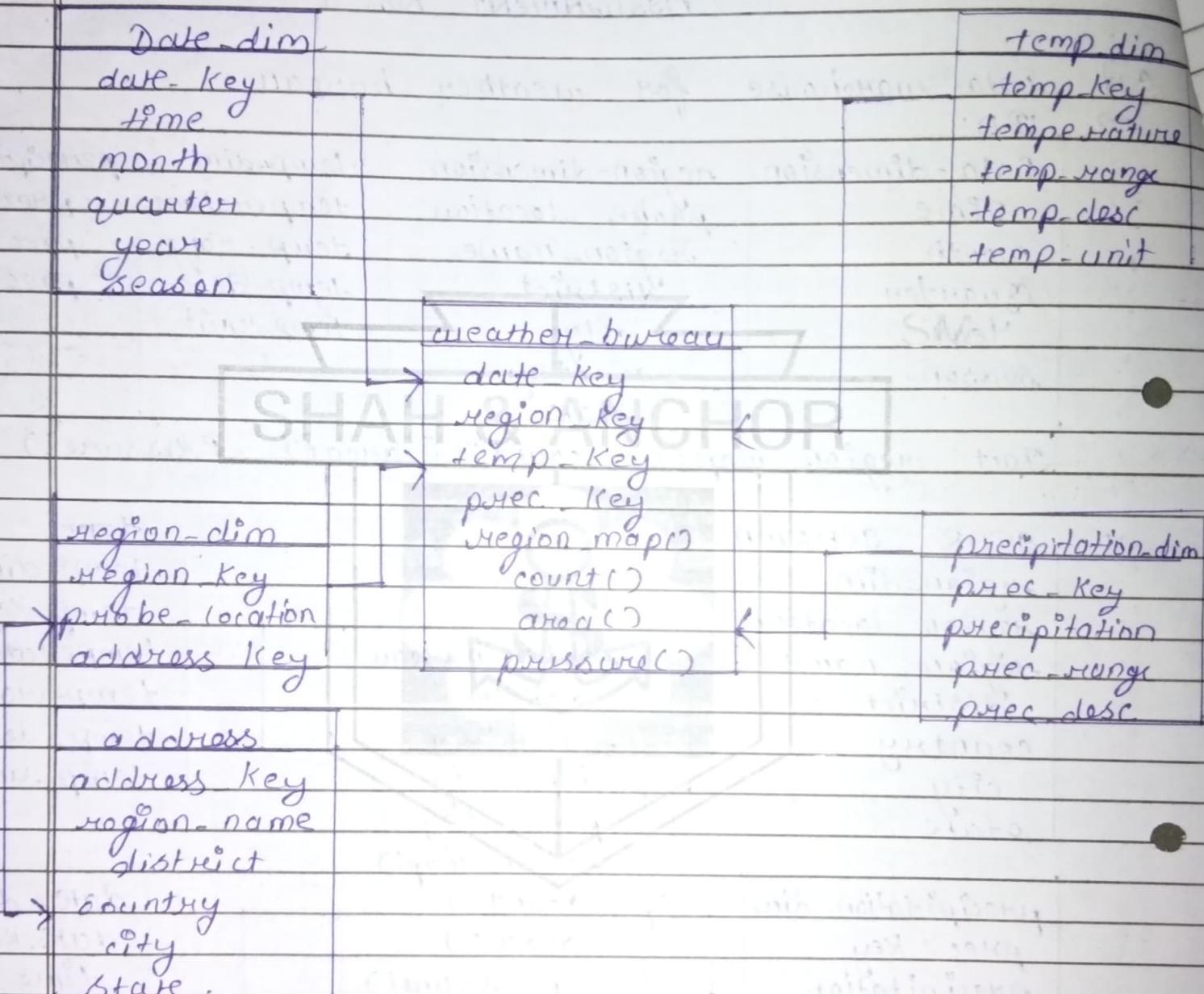
weather-bureau  
 date-key  
 region-key  
 temp-key  
 prec-key  
 region-map()

count()  
 area()  
 pressure()

temp  
 temp-dim  
 temp-key  
 temperature  
 temp-range  
 temp-desc  
 temp-unit

date-dim  
 date-key  
 time  
 month  
 quarter  
 year  
 season

## 2. SNOWFLAKE SCHEMA



Q2 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30,  
33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

a) Mean =  $\frac{809}{27} = 29.96$  years

median =  $\frac{(27+1)}{2} = 14^{\text{th}} \text{ term} = 25 \text{ years}$

b) Mode = 25. 35

Data's. modality is Bimodal

c) mid Range =  $\frac{(\text{max} + \text{min})}{2} = \frac{(70+13)}{2} = \frac{83}{2} = 41.5 \text{ year}$

d) Q1 will be median of first half of dataset  
i.e median of

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25

$\therefore Q1 = 20$  (first quartile)

Q3 will be median of second half of dataset  
i.e median of

30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52,

70

$\therefore Q3 = 35$  (third quartile)

e) 5 number summary:

min = 13

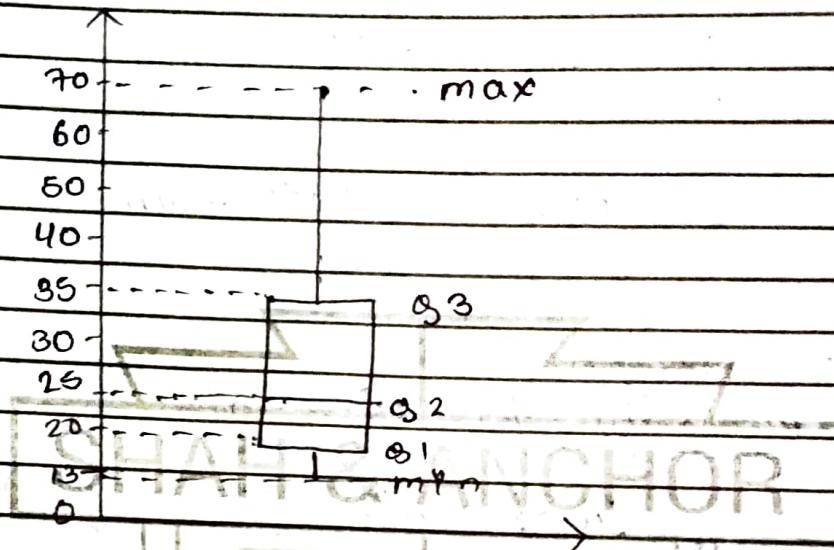
1<sup>st</sup> quartile Q1 = 20

median Q2 = 25

3<sup>rd</sup> quartile Q3 = 35

max = 70

f) Box Plot:



Q3) DEPARTMENT	AGE	SACARY	COUNT	STATUS
Sales	31...35	46K...50K	30	senior
Sales	26...30	26K...30K	40	junior
Sales	31...35	31K...35K	40	junior
systems	21...25	46K...50K	20	junior
systems	26...30	66K...70K	5	junior
systems	41...45	46K...50K	3	senior
systems	36...40	66K...70K	8	junior
marketing	31...35	46K...50K	10	senior
marketing	41...45	41K...45K	4	senior
secretary	46...50	36K...40K	4	junior
secretary	26...30	26K...30K	6	junior

use Naive Bayesian to find status of tuple:  
 ("System", "26-30", "46K.. 50K")

**MAHAVIR EDUCATION TRUST'S**  
**Shah & Anchor Kutchhi Engineering College**

Date: \_\_\_\_\_

Finding Prior Probability

$$P(c_1: \text{status} = \text{"Junior"}) \Rightarrow 6/11 = 0.55$$

$$P(c_2: \text{status} = \text{"Senior"}) = 5/11 = 0.45$$

Finding conditional probability

$$P(\text{department} = \text{"systems"} | \text{status} = \text{"Senior"}) = 2/5 = 0.4$$

$$P(\text{department} = \text{"systems"} | \text{status} = \text{"Junior"}) = 2/6 = 0.33$$

$$P(\text{age} = \text{"26... 30"} | \text{status} = \text{"Senior"}) = 0/5 = 0$$

$$P(\text{age} = \text{"26... 30"} | \text{status} = \text{"Junior"}) = 3/6 = 0.5$$

$$P(\text{Salary} = \text{"46K... 50K"} | \text{status} = \text{"Senior"}) = 2/5 = 0.4$$

$$P(\text{Salary} = \text{"46K... 50K"} | \text{status} = \text{"Junior"}) = 2/6 = 0.33$$

using above probabilities, we obtain

$$P(x | \text{status} = \text{"Junior"}) = 0.33 \times 0.5 \times 0.33 \\ = 0.05$$

$$P(x | \text{status} = \text{"Senior"}) = 0.4 \times 0 \times 0.4 = 0$$

$$P(x | \text{status} = \text{"Junior"}) \cdot P(\text{status} = \text{"Junior"}) = 0.05 \times 0.55 \\ = 0.0275 \\ = 0.03$$

$$P(x | \text{status} = \text{"Senior"}) \cdot P(\text{status} = \text{"Senior"}) = 0 \times 0.45 \\ = 0$$

$0.03 > 0$ , naive bayesian classifier predicts status of given tuple as "Junior"

Q4

Point	x	y
A <sub>1</sub>	2	10
A <sub>2</sub>	2	5
A <sub>3</sub>	8	4
B <sub>1</sub>	5	8
B <sub>2</sub>	7	5
B <sub>3</sub>	6	4
C <sub>1</sub>	1	2
C <sub>2</sub>	4	9

Initial cluster centres: A (2, 10) B (5, 8) C<sub>1</sub> (1, 2)  
 calculate the distance of each point from each cluster center

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

for A<sub>1</sub>:

$$d(A_1, A_1) = 0$$

$$d(A_1, B_1) = 3.60$$

$$d(A_1, C_1) = 8.06$$

$$d(A_3, A_1) = 8.49$$

$$d(A_3, B_1) = 5$$

$$d(A_3, C_1) = 7.28$$

$$d(B_2, A_1) = 7.07$$

$$d(B_2, B_1) = 3.60$$

$$d(B_2, C_1) = 6.70$$

$$d(C_1, A_1) = 8.06$$

$$d(C_1, B_1) = 7.21$$

$$d(C_1, C_1) = 0$$

$$d(A_2, A_1) = 5$$

$$d(A_2, B_1) = 4.24$$

$$d(A_2, C_1) = 3.16$$

$$d(B_1, B_1) = 0$$

$$d(B_1, A_1) = 3.6$$

$$d(B_1, C_1) = 7.21$$

$$d(B_3, A_1) = 7.21$$

$$d(B_3, B_1) = 4.12$$

$$d(B_3, C_1) = 5.38$$

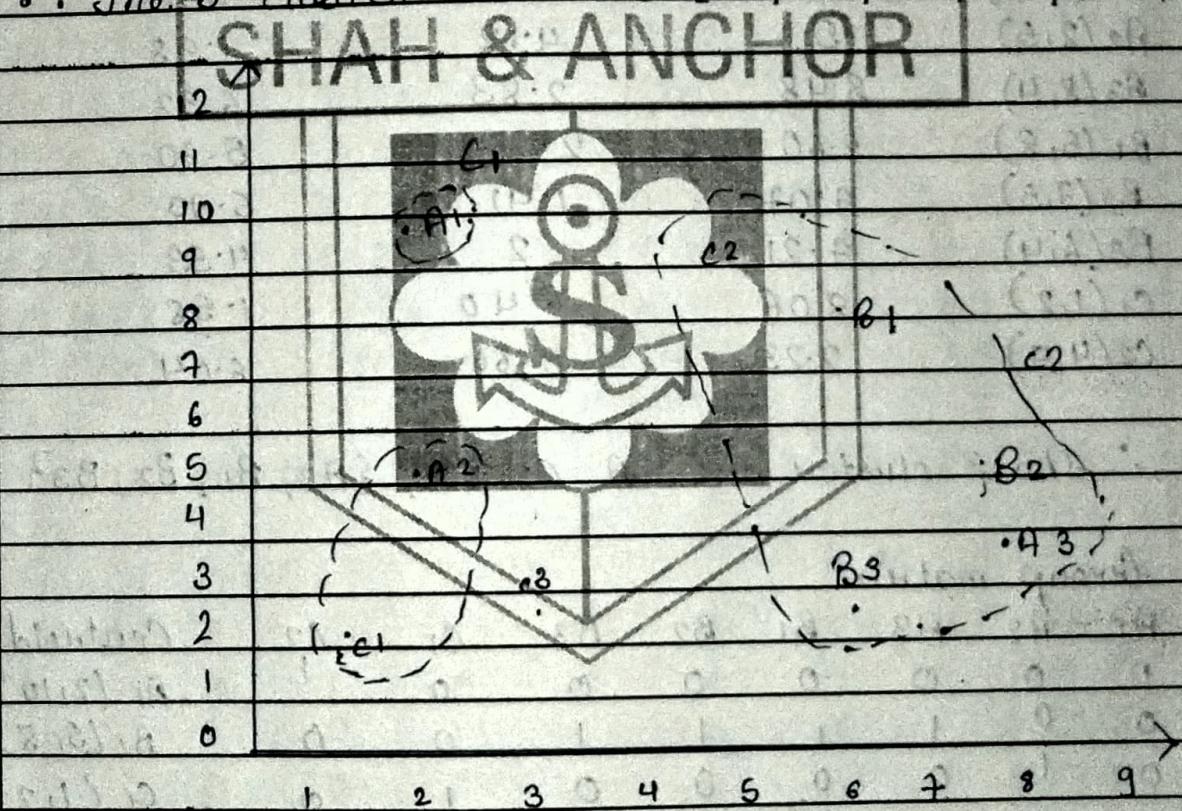
$$d(C_2, A_1) = 2.23$$

$$d(C_2, B_1) = 1.41$$

$$d(C_2, C_1) = 7.61$$

POINT	DIST. from A <sub>1</sub>	DIST. from B <sub>1</sub>	DIST. from C <sub>1</sub>	Cluster
A <sub>1</sub> (2, 10)	0	3.6	8.06	A <sub>1</sub>
A <sub>2</sub> (2, 5)	5	4.24	3.16	C <sub>1</sub>
A <sub>3</sub> (8, 4)	8.49	5	7.28	B <sub>1</sub>
B <sub>1</sub> (5, 8)	3.6	0	7.21	B <sub>1</sub>
B <sub>2</sub> (7, 5)	7.07	3.6	6.7	B <sub>1</sub>
B <sub>3</sub> (6, 4)	7.21	4.12	5.38	B <sub>1</sub>
C <sub>1</sub> (1, 2)	8.06	7.21	0	C <sub>1</sub>
C <sub>2</sub> (4, 9)	2.23	1.41	7.61	B <sub>1</sub>

∴ The 3 cluster are: {A<sub>2</sub>}, {A<sub>3</sub>, B<sub>1</sub>, B<sub>2</sub>, B<sub>3</sub>, C<sub>2</sub>}, {A<sub>1</sub>, C<sub>1</sub>}



Group Matrix a:

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	C <sub>1</sub>	C <sub>2</sub>	Centroid
1	0	0	0	0	0	0	0	A <sub>1</sub> (2, 10)
0	0	1	1	1	1	0	1	B <sub>1</sub> (5, 8)
0	1	0	0	0	0	1	0	C <sub>1</sub> (1, 2)

∴ new centroids after first iteration

$$C_1 = (2, 10)$$

$$C_2 = \left( \frac{8+5+7+6+4}{5}, \frac{4+8+5+4+9}{5} \right) = (6, 6)$$

$$C_3 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

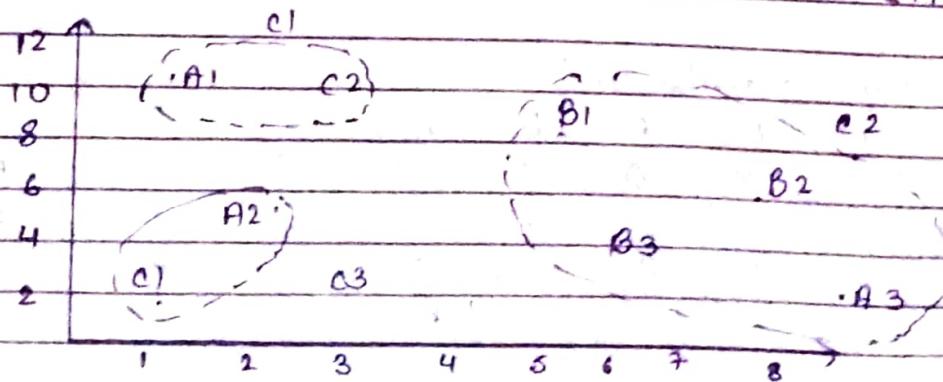
Iteration ~ 2:

Point	Dist from C1	Dist from C2	Dist from C3	Cluster
A <sub>1</sub> (2, 10)	0	5.65	6.62	C <sub>1</sub>
A <sub>2</sub> (2, 5)	5	4.12	1.58	C <sub>3</sub>
A <sub>3</sub> (8, 4)	8.48	2.83	6.52	C <sub>2</sub>
B <sub>1</sub> (5, 8)	3.60	2.23	5.70	C <sub>2</sub>
B <sub>2</sub> (7, 5)	7.07	1.41	5.70	C <sub>2</sub>
B <sub>3</sub> (6, 4)	7.21	2	4.53	C <sub>2</sub>
C <sub>1</sub> (1, 2)	8.06	6.40	1.58	C <sub>3</sub>
C <sub>2</sub> (4, 9)	2.23	8.66	6.04	C <sub>1</sub>

∴ The 3 clusters are {A<sub>1</sub>, C<sub>2</sub>}, {A<sub>3</sub>, B<sub>1</sub>, B<sub>2</sub>, B<sub>3</sub>} {A<sub>2</sub>, C<sub>1</sub>}

Group matrix:

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	C <sub>1</sub>	C <sub>2</sub>	Centroid
1	0	0	0	0	0	0	-1	A <sub>1</sub> (2, 10)
0	0	1	1	1	1	0	0	B <sub>1</sub> (5, 8)
0	1	0	0	0	0	1	0	C <sub>1</sub> (1, 2)



MAHAVIR EDUCATION TRUST'S  
Shah & Anchor Kutchhi Engineering College

Date:

N TRUST'S  
Engineering

∴ new centroids after 2nd iteration

$$C_1 = \left( \frac{2+4}{2}, \frac{10+9}{2} \right) = (3, 9.5)$$

$$C_2 = \left( \frac{8+5+7+6}{4}, \frac{4+8+5+4}{4} \right) = (6.5, 5.25)$$

$$C_3 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

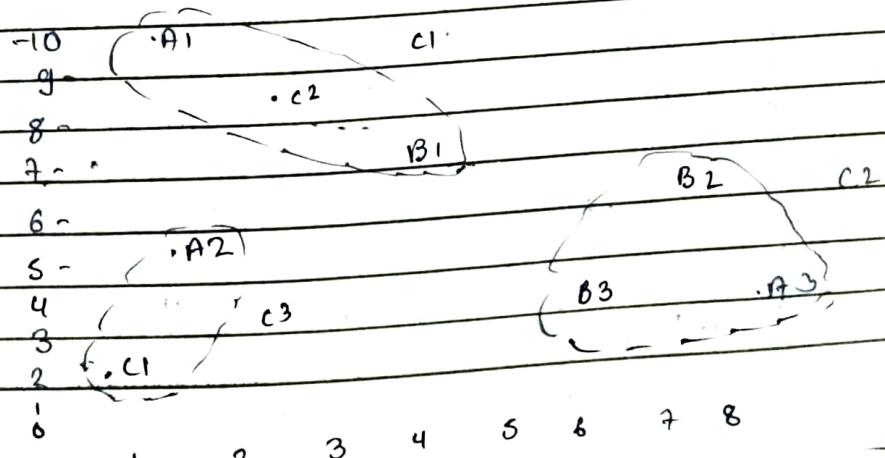
Iteration 3:

POINT	Dist from C1	Dist from C2	Dist from C3	cluster
A <sub>1</sub> (2, 10)	1.12	6.54	6.52	C <sub>1</sub>
A <sub>2</sub> (2, 5)	4.61	4.51	1.58	C <sub>3</sub>
A <sub>3</sub> (8, 4)	7.43	1.95	6.52	C <sub>2</sub>
B <sub>1</sub> (5, 8)	2.5	3.13	5.70	C <sub>1</sub>
B <sub>2</sub> (7, 5)	6.02	0.56	5.70	C <sub>2</sub>
B <sub>3</sub> (6, 4)	6.26	1.35	4.53	C <sub>3</sub>
C <sub>1</sub> (1, 2)	7.76	6.39	1.58	C <sub>3</sub>
C <sub>2</sub> (4, 9)	1.12	4.61	6.04	C <sub>1</sub>

∴ The cluster are {A<sub>1</sub>, B<sub>1</sub>, C<sub>2</sub>}, {A<sub>3</sub>, B<sub>2</sub>, B<sub>3</sub>}, {A<sub>2</sub>, C<sub>1</sub>}

Group matrix

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	C <sub>1</sub>	C <sub>2</sub>	Centroid
1	0	0	1	0	0	0	0	1	A <sub>1</sub> (2, 10)
0	0	1	0	1	1	1	0	0	B <sub>1</sub> (5, 8)
0	1	0	0	0	0	0	1	0	C <sub>1</sub> (1, 2)



∴ new centroids after 3rd iteration are:

$$C_1 = \left( \frac{2+5+4}{3}, \frac{9+8+10}{3} \right) = (3.67, 9)$$

$$C_2 = \left( \frac{8+7+6}{3}, \frac{4+5+4}{3} \right) = (7, 4.33)$$

$$C_3 = (1.5, 3.5)$$

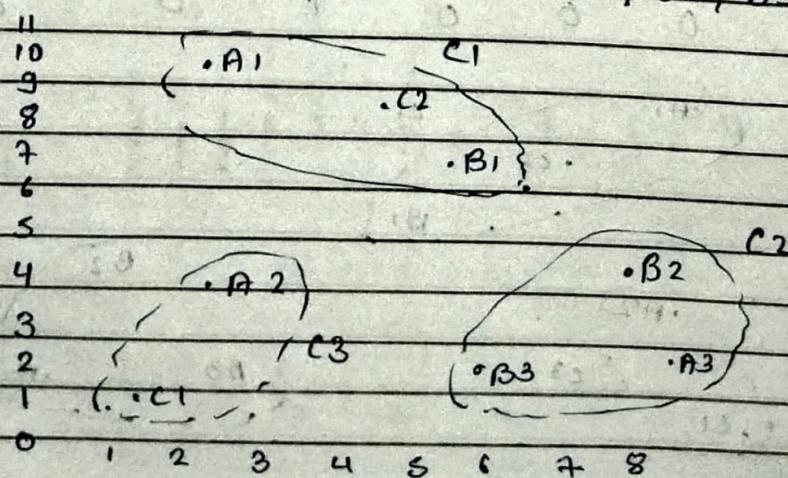
Point	Dist from C1	Dist from C2	Dist from C3	Cluster
A <sub>1</sub> (2,10)	1.95	7.56	6.52	C <sub>1</sub>
A <sub>2</sub> (2,5)	4.33	5.04	1.58	C <sub>3</sub>
A <sub>3</sub> (8,4)	6.61	1.05	6.52	C <sub>2</sub>
B <sub>1</sub> (5,8)	1.66	4.18	5.10	C <sub>1</sub>
B <sub>2</sub> (9,5)	5.20	0.67	5.70	C <sub>2</sub>
B <sub>3</sub> (6,4)	5.51	1.05	4.53	C <sub>2</sub>
C <sub>1</sub> (1,2)	7.49	6.44	1.58	C <sub>3</sub>
C <sub>2</sub> (4,9)	0.64	5.55	6.04	C <sub>1</sub>

The 3 clusters are: {A<sub>1</sub>, B<sub>1</sub>, C<sub>2</sub>} {A<sub>3</sub>, B<sub>2</sub>, B<sub>3</sub>} {A<sub>2</sub>, C<sub>1</sub>}

Group matrix:

A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	C <sub>1</sub>	C <sub>2</sub>	Centroid
1	0							
0	0							
0	1							

cluster in iteration 4 is same as in iteration 3  
 $\therefore$  Final cluster {A<sub>1</sub>, B<sub>1</sub>, C<sub>2</sub>} {A<sub>3</sub>, B<sub>2</sub>, B<sub>3</sub>} {A<sub>2</sub>, C<sub>1</sub>}



85

## APRIORI

TID	Items bought
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, F, Y}
T300	{M, A, K, E}
T400	{M, U, C, K, Y}
T500	{C, D, O, K, I, E}

minimum support = 60%, minimum confidence = 80%.

Support = Juples containing (A+B)

no. of tuples

60  
100

$\geq$   
5

$x = 3$  i.e. min support count = 3

item	Count
m	3
o	4
n	2
k	5
g	4
y	3
d	1
a	1
u	1
c	2
i	1

$C_1 =$

g  
y  
d  
a  
u  
c  
i

After pruning, we get: (count  $\geq 3$ )

Item	Count
L1 = m	3
n	3
K	5
E	4
Y	3

Generating C2 from L1:

Item	Count
MO	1
MK	3
ME	2
L1 $\bowtie$ L1	
mY	2
OK	3
OE	3
04	2
KE	4
KY	3
EV	2

After pruning, we get (count  $\geq 3$ )

Item	Count
L2	
$\rightarrow$	
MK	3
OK	3
OE	3
KE	4
KY	3

Generate  $C_3$  from  $L_2$ :

Item	Count
OKE	3
KEY	2
OEY	2
MKY	2
MOE	1
MKO	1
KOY	2

After pruning, we get (Count  $\geq 3$ )

Item	Count
OKE	3

$\therefore$  frequent item set =  $\{O, K, E\}$

$\therefore$  we generate association rules for  $\{O, K, E\}$

Confidence =  $\frac{\text{Juples containing } A \text{ & } B}{\text{Juples containing only } A}$

$\therefore$  association rules are:

	Confidence
$O \wedge K \Rightarrow E$	$3/3 = 100\%$
$O \wedge E \Rightarrow K$	$3/3 = 100\%$
$E \wedge K \Rightarrow O$	$3/4 = 75\%$
$E \wedge O \Rightarrow O \wedge K$	$3/4 = 75\%$
$K \Rightarrow O \wedge E$	$3/5 = 60\%$
$O \Rightarrow E \wedge K$	$3/3 = 100\%$

Now, we have to find strong association rule i.e  
confidence = 80%.

∴ The strong association rules are:

$$O \wedge K \Rightarrow E$$

$$O \wedge E \Rightarrow K$$

$$O \Rightarrow E \wedge K$$

Q6 Describe HITS algorithm

- A clever aimed at finding both authoritative pages & hubs
  - Hub is a page that contains links to authoritative pages
  - The clever system identifies authoritative pages & hubs by creating weights
  - HITS stand for Hyperlink-Induced topic search
  - It is used to the web-link structures to discover & rank the webpage relevant for a particular search
1. Based on a given set of keywords, a set of relevant pages is found
  2. Hub & authority measures are associated with these pages. Pages with highest values are returned.
- $w \sqcup www$  viewed as directed graph

$g \sqcup$  query

$s \sqcup$  support

OUTPUT:

$A \sqcup$  set of authority pages

$H \sqcup$  set of hub pages

HITS algorithm

$R = SE(w, q)$

$B = R \cup \{ \text{pages linked to from } R \} \cup \{ \text{pages that link to page in } R \}$

$G = (B, L) =$  subgraph of  $w$  induced by  $B$

$G' = (B, L) =$  Delete links in  $G$  within the same site

$XP = g$  where  $(g, \text{plot}) \uparrow q$ ;  $\sqcup$  find authority of weights

$Y_{pq} = \text{where } (p, a) \text{ } p \in \text{ng}; \sqcup$  find hub weights

$A = \{ p \mid p \text{ has one the height } xp \};$

$H = \{ p \mid p \text{ has one of the height } up \};$