

NAME: SHIV PRASAD · C. PREMARAJAN

ROLL NO: 9696

SUB: DWM

ELONIC COLLEGE OF ENGINEERING

DATE
26/10/23

Assignment 1

Q.1] A bank wants to develop a DW for effective decision making about their loan schemes. Loans are provided for various purpose (Home, car, Education, etc).

The whole country is categorized into regions namely East, west, South, North. Loan is distributed to customers at the interest rate that changes from time to time. Different loans have different ROI. Clearly design the star schema & snowflake schema.

Information Package Diagram.

Business Process: Loan Analyst.

Time	Customer	Region	Loan
Trans-key	Customer-key	Region-key	loankey
Day	Account-number	Region-name	loan-type
Day-of-month	Account-type		ROI
month	Customer-name		Payment
Quarter			
Year			

Fact: Loan-amount, No.of: customer.

Star Schema

Time Dimension table

Time - Key
Day
Day-of-week
Month
Quarter
Year

Region Dimension table

Region - Key
Region - name

Fact table

Time - Key
Customer - Key
Region - Key
Loan - Key
Loan - amount
No. of Customer

Customer

Customer -
Customer -
Account - num
Account - type

Loan Dimension table

Loan - Key
Loan - type
ROI
Payment

* Snowflake Schema

Time Dimension Table

Time - Key
Day
Day-of-week
Month
Quarter
Year

Region Dimension Table

Region - Key
Region - name
Branch Key

Fact table

Time - Key
Customer - key
Region - key
Loan - Key
Loan - amount
No. of Customer

Customer Dimension table

Customer - Key
Account - number
Account - type

Loan Dimension table

Loan - Key
ROI
Payment - Key

Branch Dimension Table

Branch - Key
Branch - name
City
State

Payment Dimension table

Payment - Key
RBI Banking
Diamond Delft
NBFT

2] Consider a data warehouse for weather related data like region (state, city, area) time (year, quarter, month) and temperature (temp - category, weather-condition, season) with facts as avg-unit-int and avg humidity. Using this example perform following.

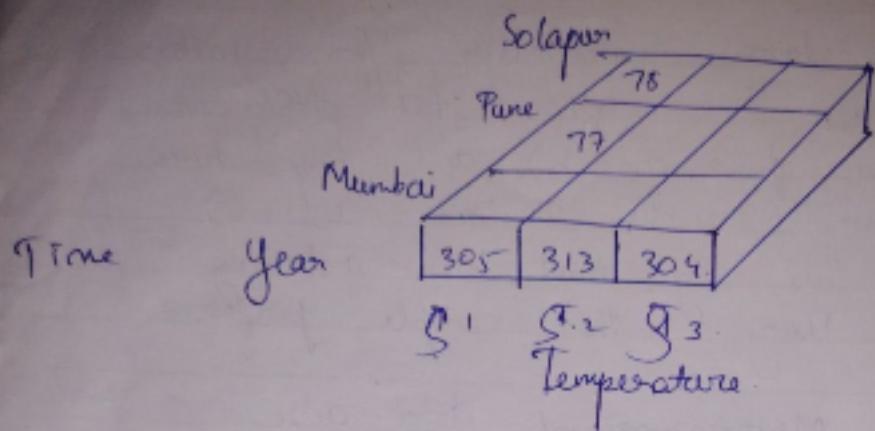
- 1) Draw a Multidimensional data cube.
- 2) perform following OLAP operations:
 - Roll-up.
 - Drill-down.
 - Slice
 - Dice
 - Pivot.

		Solapur	78		
		Pune	77		
		Mumbai			
Time	Q1	75	78	70	
Dimension	Q2	85	89	84	
	Q3	78	81	80	
	Q4	62	68	70	
		Sum	32	33	

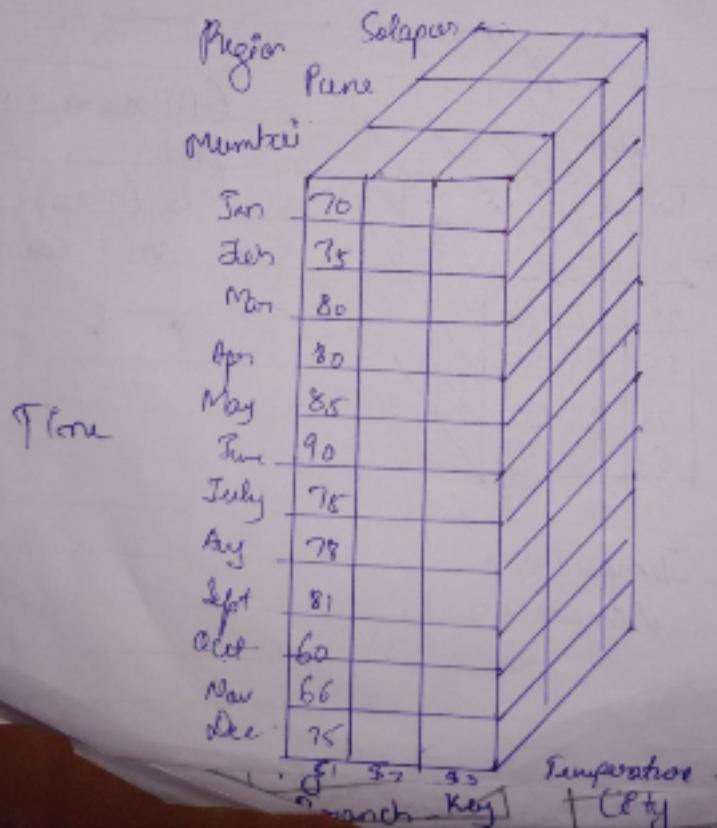
(Measure i.e.
fact
is Average
visit count).

Temperature
Dimension

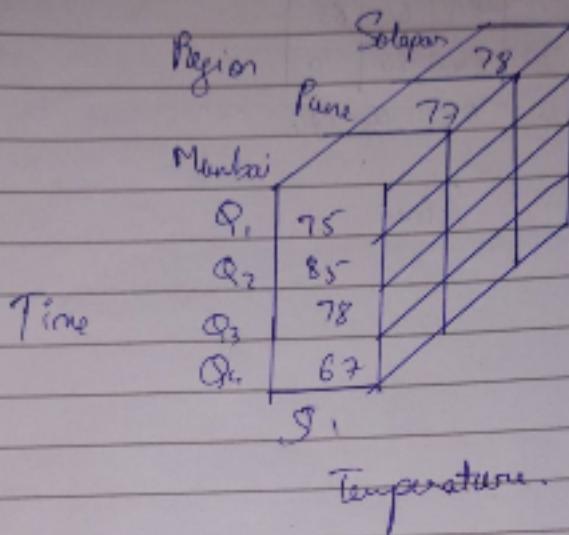
1) Roll-up (Time)



1) Drill-down (Time)



S196 (Temperature = T_1)

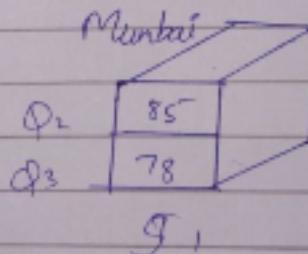


Dice
=

Time = Q₂ Q₃ Q₃

Temperature = T_1

Region = Mumbai



Pivot

=

Temperature

\bar{S}_1	75	85	78	67
\bar{S}_2	78	87	81	65
\bar{S}_3	70	84	80	70
	Q_1	Q_2	Q_3	Q_4

Time

~~MAPLE~~
BWM Assignment no - 2

Q1] The following Data (in increasing order) for the attribute age : 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.

→

Partition into (equal-depth) bins.

Bin1 :- 13, 15, 16, 16, 19, 20, 20, 21, 22.

Bin2 :- 22, 25, 25, 25, 25, 30, 33, 33, 35.

Bin3:- 35, 35, 35, 36, 40, 45, 46, 52, 70.

a) Smoothing by bin mean-

Bin1:- 18, 18, 18, 18, 18, 18, 18, 18, 18.

Bin2:- 28, 28, 28, 28, 28, 28, 28, 28, 28.

Bin3:- 44, 44, 44, 44, 44, 44, 44, 44, 44.

b) Smoothing by bin median-

Bin1:- 19, 19, 19, 19, 19, 19, 19, 19, 19.

Bin2:- 25, 25, 25, 25, 25, 25, 25, 25, 25.

Bin3:- 40, 40, 40, 40, 40, 40, 40, 40, 40.

Q.2] Explain with examples Sampling methods:-

a) Simple random sample without replacement

Simple random sampling without replacement (SRSWOR)
is a sampling method in which each unit in the population has an equal chance of being selected and each unit can only be selected once.

E:- Suppose in class of 20 students and you want to select a sample of 5 students to participate in debate.
We can use SRSWOR to select the sample.

Step1:- Write each student's name on slip of paper.

Step2:- put the slips of paper in a hat or bowl & mix them.

Step3:- Draw 5 slips of paper without replacement.

Step4:- The students whose name are drawn will be the sample.

b) Simple random sample with replacement

→ Simple random sampling with replacement (SRSWR)
is a sampling method in which each unit in the population has an equal chance of being selected and each unit can be selected multiple times.

E:- Suppose we have a bag of 10 marbles, 5 red & 5 blue. You want to draw a sample of 3 marbles with replacement.

Step1:- Randomly select one marble from bag.

Step2:- Record the color of marble & put it back in bag.

Step3:- Repeat steps 1 & 2 until you have a sample of 3 marbles
P/O:- Red, Blue, Red.

c) Cluster Sample.

→ Cluster sampling is a sampling method in which population is divided into groups or clusters. A random sample of cluster is selected. All units within the selected clusters are then included in the sample.
e.g:-

Suppose we want to survey the opinions of students at a university. We could divide the university into the clusters by college. So then randomly select a sample of colleges & survey all the students within these colleges.

d) Stratified Sample.

→ A Stratified Sample is a sample that is divided into groups or strata based on specific characteristics. The characteristics used to stratify the sample is called the stratification variable. Once the sample is stratified, a random sample is taken from each stratum.
e.g:-

Suppose you want to survey the opinions of student at a university. stratify by gender (M/F). Then randomly select the sample of Male & female from each college at university.

Q.3] Normalise the following group of data.

→ 200, 300, 400, 500, 1000.

→ a) Min-max normalization by setting min 0 & max 1.

→ min = 200, max = 1000, new-min = 0, new-max = 1.

$$X = \frac{x - \text{min}}{\text{max} - \text{min}} (1 - 0) + 0 = (x - \text{min})$$

$$x_1 = \frac{200 - 200}{1000 - 200} (1 - 0) + 0 = 0$$

$$x_2 = \frac{300 - 200}{1000 - 200} (1 - 0) = 0.125$$

$$x_3 = \frac{400 - 200}{1000 - 200} (1 - 0) + 0 = 0.25$$

$$x_4 = \frac{500 - 200}{1000 - 200} (1 - 0) + 0 = 0.5$$

$$x_5 = \frac{1000 - 200}{1000 - 200} (1 - 0) + 0 = 1$$

Normalised values are 0, 0.125, 0.25, 0.5, 1.

b) Z-score Normalization.

$$Z = \frac{x - \mu}{\sigma}$$

$$\bar{\mu} = \frac{200 + 300 + 400 + 600 + 1000}{5} = 500$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{\mu})^2}{N}}$$

$$\sigma = \sqrt{\frac{(200-500)^2 + (300-500)^2 + (400-500)^2 + (600-500)^2 + (1000-500)^2}{5}}$$

$$\sigma = 282.84$$

$$Z_1 = \frac{x_1 - \mu}{\sigma} = \frac{200 - 500}{282.84} = -1.0606$$

$$Z_2 = \frac{x_2 - \mu}{\sigma} = \frac{300 - 500}{282.84} = -0.7071$$

$$Z_3 = \frac{x_3 - \mu}{\sigma} = \frac{400 - 500}{282.84} = -0.3538$$

$$Z_4 = \frac{x_4 - \mu}{\sigma} = \frac{600 - 500}{282.84} = 0.3535$$

$$Z_5 = \frac{x_5 - \mu}{\sigma} = \frac{1000 - 500}{282.84} = 1.7677$$

C) Normalization by decimal scaling.

$n' = \frac{n}{10^j}$ where j is the smallest integer such that $\max(|v'|) < 1$.

10^4 is 10000. ($j=4$ as we have maximum 4 digit number in the range).

$$n'_1 = \frac{200}{10^4} = \frac{200}{10000} = 0.02.$$

$$n'_2 = \frac{300}{10^4} = \frac{300}{10000} = 0.03$$

$$n'_3 = \frac{400}{10^4} = \frac{400}{10000} = 0.04.$$

$$n'_4 = \frac{600}{10^4} = \frac{600}{10000} = 0.06.$$

$$n'_5 = \frac{1000}{10^4} = \frac{1000}{10000} = 0.1.$$

Normalized values are 0.02, 0.03, 0.04, 0.06, 0.1.

Q. 4.

a) Min-Max normalization to transform the value 35 for age onto the range [0.0, 1.0]

$$\hookrightarrow n' = \frac{n - \text{min}}{\text{max} - \text{min}} \quad [\text{new-min} - \text{new-max}] + [\text{new-min}]$$

$$\text{Min} = 13$$

$$\text{Max} = 70, \text{ new-min} = 0.0, \text{ new-Max} = 1.0$$

$$n' = \frac{35 - 13}{70 - 13} [1.0 - 0.0] + [0.0]$$

$$x' = 0.3894.$$

b]. Use z-score normalization to transform the value 35 for age, where standard deviation of age is 12.94 yrs.

\rightarrow

$$z = \frac{x - \mu}{\sigma}, \mu = 29.96, \sigma = 12.94.$$

$$\therefore z = \frac{35 - 29.96}{12.94} = 0.3894.$$

c]. Use normalization by decimal scaling to transform the value 35 for age.

$$\rightarrow n' = \frac{x}{10^j}$$

10^j is 100 ($j=2$ we have maximum 2 digit number in range).

$$\therefore n' = \frac{35}{10^2} = 0.35$$

d] Min-max normalization would be done in the value range of 0 to 1, reducing the outliers. Hence it's the better technique.

Q.5]

a] Stepwise forward selection.

- i) Start with empty set of attributes.
- ii) Determine the best of original attributes and add it to set.
- iii) At each step, find the best of remaining original attributes and add it to set.

∴ Initial attribute set -

$$\{A_1, A_2, A_3, A_4, A_5, A_6\}.$$

Initial reduced set: $\{ \}$

⇒ $\{A_1\}$ - (most imp.)

⇒ $\{A_1, A_2\}$

⇒ Reduced attribute set $\{A_1, A_2, A_3\}$

b] Stepwise backward selection.

∴ i) Starts with full set of attributes.

ii) At each step, it removes the worst attribute (having least imp.) from set.

∴ Initial attribute set -

$$\{A_1, A_2, A_3, A_4, A_5, A_6\}$$

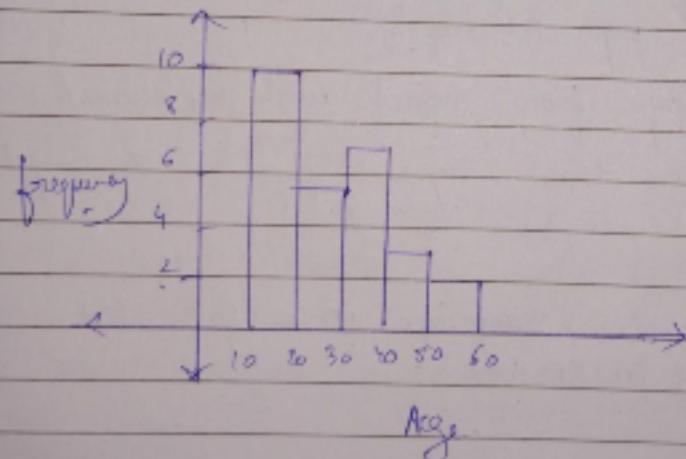
Initial reduced set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$ - {~~Choosing least imp.
A₁ & A₂ att.~~}

⇒ Reduced attribute set: $\{A_3, A_4, A_5\}$.

- c) A combination of forward selection & backward selection.
- i) The procedure combining & selects the best attributes & removes the worst among the binary attributes.
- ii) For all above method stopping criteria is different and it requires a threshold on the measure used to stop the attribute selection process.

Q.6]

- a) Plot an equal-width histogram of width 10



b) Sketch examples of each of the following sampling techniques : SRSWOR, SRSWR.

- a) Simple Random Sampling without replacement (SRSWR)
- 1) Write each of the age attribute on a separate slip of paper.
 - 2) Put all the slips of paper in bowl.
 - 3) Mix the slips of paper.
 - 4) Draw 5 slips of paper from bowl.
 - 5) The age attribute on drawn slips of paper.
 $P/O = 20, 22, 25, 35, 46$.

b) Simple Random Sampling with Replacement (SRSWR).

- 1) Write each of the age attribute on a separate slip of paper.
- 2) Put all the slips of paper in bowl.
- 3) Mix the slips of paper thoroughly.
- 4) Draw 5 slips of paper from bowl, with replacement.
- 5) The age attribute on drawn slips of paper on the sample.
 $P/O = 25, 25, 35, 20, 35$.

Note:- possible to draw same age in SRSWR but
not in SRSWOR.

c) Cluster Sampling & Stratified Sampling - use same of size 5 and the strata "youth", "middle aged" and "senior".

→ a) Cluster Sampling.

- 1) Cluster 1: Youth (age < 25)
- 2) Cluster 2: Middle - age (age ≥ 25 & age < 65)
- 3) Cluster 3: Senior (age ≥ 65)

Next we can randomly select one cluster from each of the three strata. Finally we can select all the age attributes from selected cluster to form our sample.

$$\begin{array}{ll} p/0: & \text{Cluster 1: 19} \\ p/2: & \text{Cluster 2: 30, 35} \\ & \text{Cluster 3: 70.} \end{array}$$

b) Stratified sampling.

- 1) Cluster 1: Youth (age < 25)
- 2) Cluster 2: middle - age (age ≥ 25 and age < 65)
- 3) Cluster 3: Senior (age ≥ 65)

Next, we can randomly select 5 age attributes from each stratum.

$$\begin{array}{l} \text{Cluster 1: 15, 19} \\ \text{Cluster 2: 25, 30, 35} \\ \text{Cluster 3: 70.} \end{array}$$

Note:

2

Cluster sampling & stratified sampling are two different sampling methods. Cluster sampling is a method of selecting groups of individual while Stratified Sampling is a method of selecting individual from different groups.

Name: Chavaretil Shivprasad Barve
Roll no: 9696

F.R. CONCESSIONS COLLEGE OF ENGINEERING

Assignment 3:

Varun

Q.1] Apply Naive Bayesian Classification Algorithm on following data

→ tuple (Homeowner = Yes, Status = Employed, Income = Avg)

$$p(\text{Defaulted} = \text{'yes'}) = 3/10$$

$$p(\text{Defaulted} = \text{'No'}) = 7/10$$

To calculate probability

$$P(C_i|x) = \frac{p(x|C_i)p(C_i)}{p(x)}$$

∴ $p(x)$ is constant for all classes

$$\therefore p(C_i|x) = p(x|C_i)p(C_i)$$

$$p(x|C_i) = p(x|\text{Defaulted} = \text{'yes'})$$

$$= p(\text{Homeown} = \text{'Yes'} | \text{Defaulted} = \text{'yes'}) \times$$

$$= p(\text{Status} = \text{'Employed'} | \text{Defaulted} = \text{'yes'}) \times$$

$$= p(\text{Income} = \text{'Average'} | \text{Defaulted} = \text{'yes'})$$

$$= 0 \times \frac{2}{3} \times 0 = 0$$

$$= 0$$

$$p(\text{Defaulted} = \text{'yes'} | x) =$$

$$p(x| \text{Defaulted} = \text{'yes'}) \times p(\text{Defaulted} = \text{'yes'})$$

$$= 0 \times \frac{3}{10} = 0.$$

$$\begin{aligned}
 p(x|c_i) &= p(x | \text{Defaulted} = 'No') \\
 &= p(\text{Homeowner} = 'Yes' | \text{Defaulted} = 'No') \times \\
 &\quad p(\text{Status} = 'Employed' | \text{Defaulted} = 'No') \times \\
 &\quad p(\text{Income} = 'Average' | \text{Defaulted} = 'No') \\
 &= \frac{2}{7} \times \frac{2}{7} \times \frac{1}{7} = \frac{6}{7^3}
 \end{aligned}$$

$$\begin{aligned}
 \therefore p(\text{Defaulted} = 'No' | x) &= \\
 &= p(x | \text{Defaulted} = 'No') \times p(\text{Defaulted} = 'No') \\
 &= \frac{6}{7^3} \times \frac{7}{10} = \frac{6}{50} \\
 &= 0.012
 \end{aligned}$$

$$\therefore p(\text{Defaulted} = 'No' | x) > p(\text{Defaulted} = 'Yes' | x)$$

∴ Therefore the naive bayes classifier predicts
 Defaulted = 'No' for sample x.

Q.2. Apply naive bayesian classification algorithm on following data tuple ($\text{District} = \text{'Rural'}$, House type = 'Semi-detached' , Income = 'Low' , Previous customer = 'No').

$$P(\text{Outcome} = \text{'Nothing'}) = \frac{5}{14}$$

$$P(\text{Outcome} = \text{'Reponded'}) = \frac{9}{14}$$

To calculate probability,

$$P(c_i|x) = \frac{p(x|c_i)p(i)}{p(x)}$$

$\therefore p(x)$ is constant for all classes

$$\therefore P(c_i|x) = p(x|c_i)p(c_i)$$

$$\begin{aligned} \therefore P(x|c_i) &= P(x | \text{Outcome} = \text{'Nothing'}) \\ &= P(\text{District} = \text{'Rural'} | \text{Outcome} = \text{'Nothing'}) \times \\ &\quad P(\text{House Type} = \text{'Semi-detached'} | \text{Outcome} = \text{'Nothing'}) \\ &= P(\text{Income} = \text{'Low'} | \text{Outcome} = \text{'Nothing'}) \times \\ &= P(\text{Previous customer} = \text{'No'} | \text{Outcome} = \text{'Nothing'}) \\ &= \frac{0}{5} \times \frac{1}{5} \times \frac{2}{5} \end{aligned}$$

$$\therefore p(\text{Outcome} = \text{'Nothing'} | x) = p(x | \text{Outcome} = \text{'Nothing'}) \\ = x p(\text{Outcome} = \text{'Nothing'}) \\ = 0 \times \frac{5}{14} = 0$$

$$\therefore p(x | c_i) = p(x | \text{Outcome} = \text{'Responded'}) \\ = p(\text{District} = \text{'Rural'} | \text{Outcome} = \text{'Responded'})x \\ = p(\text{House Type} = \text{'Semi-detached'} | \text{Outcome} = \text{'Responded'}) \\ = p(\text{Income} = \text{'low'} | \text{Outcome} = \text{'Responded'})x \\ = p(\text{Previous Customer} = \text{'No'} | \text{Outcome} = \text{'Responded'}) \\ = \frac{4}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9} \\ = \frac{64}{729}$$

$$p(\text{Outcome} = \text{'Responded'} | x) = p(x | \text{Outcome} = \text{'Responded'}) \\ \times p(\text{Outcome} = \text{'Responded'}) \\ = \frac{64}{729} \times \frac{9}{14} = 0.056.$$

$\therefore p(\text{Outcome} = \text{'Responded'} | x) > p(\text{Outcome} = \text{'Nothing'} | x)$

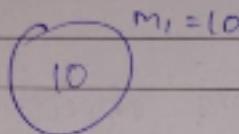
Therefore the naive bayes classifier predicted
 $\text{Outcome} = \text{'Responded'}$ for our sample x .

Q.3] Cluster the following data (10, 16, 6, 29, 16, 38, 42, 8, 11).

Initial K=2

→ Find centroids bco $K_1 = 10 \text{ & } K_2 = 29$.

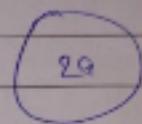
Iteration 1 K_1



$$m_1 = 10$$

$$K_1 = \{10, 18, 6, 16, 8, 11\}$$

K_2



$$m_2 = 29$$

$$K_2 = \{29, 38, 42\}$$

Update the centroids.

$$\therefore C_1 = \frac{10 + 15 + 6 + 16 + 8 + 11}{6} \\ = 11$$

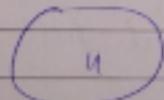
$$C_2 = \frac{29 + 38 + 42}{3} \\ = 36.33.$$

\therefore New centroids & assumed centroids are not equal,
 \therefore therefore repeat the step

Iteration 2

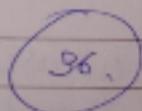
K_1

$$m_1 = 11$$



K_2

$$m_2 = 36$$



update the centroids

$$\therefore C_1 = \frac{10+15+6+16+8+11}{6}$$

$$\approx 11.$$

$$C_2 = \frac{24+38+42}{3} \\ = 36.$$

- \therefore The centroids did not change in this iteration
 \therefore The clusters are

$$K_1 = \{10, 15, 6, 16, 8, 11\}$$

$$K_2 = \{24, 38, 42\}$$

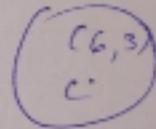
Q.4] Cluster the following data objects A(1,2) B(3,3)
C(6,3) D(2,3) E(4,4) F(1,5) where k=2.

\rightarrow Let centroids be $C_1 = A(1,2)$, $C_2 = C(6,3)$.

K₁



K₂



$$K_1 = \{A, B, D, F\}$$

update the centroid

$$C_1 = \left\{ \frac{1+3+2+1}{4}, \frac{2+4+3+5}{4} \right\}$$

$$= (1.75, 3.5)$$

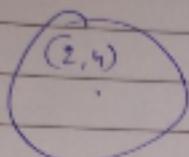
$$K_2 = \{C, E\}$$

$$C_2 = \left(\frac{6+3}{2}, \frac{3+4}{2} \right) \\ = (4.5, 3.5)$$

\therefore New centroids & assumed centroids are not same, therefore repeat the step.

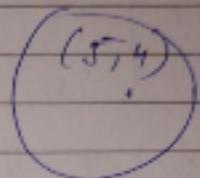
Iteration 2:

k_1



$$K_1 = \{A, B, D, F\}$$

k_2



$$K_2 = \{C, E\}$$

update the centroid

$$C_1 = \left(\frac{1+3+2+1}{4}, \frac{2+4+3+5}{4} \right)$$

$$C_1 = (1.75, 3.5) \approx (2, 4)$$

$$C_2 = \left(\frac{6+4}{2}, \frac{3+4}{2} \right)$$

$$= (5, 3.5) \approx (5, 4)$$

\therefore The centroids did not change in the iteration

\therefore The clusters are

$$K_1 = \{A, B, D, F\}$$

$$K_2 = \{C, E\}$$

Q.5) Cluster the following data using agglomeration hierarchical clustering technique

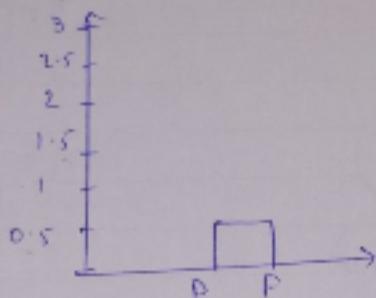
	x_1	x_2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

1) Single linkage

→ Distance matrix

	A	B	C	D	E	F
A	0					
D		0				
B	1	0				
C	8	7	0			
D	5	4	3	0		
E	6	5	2	1	0	
F	4.5	3.5	3.5	0.5	1.5	0

Point F & D has minimum distance 0.5.



Pre-computing the distance matrix $\text{dist}(D, F) \approx 0.5$

$$\text{dist}((D, F), A)$$

$$= \min(\text{dist}(D, A), \text{dist}(F, A))$$

$$= \min(5, 4.5)$$

$$\approx 4.5$$

$$\text{dist}((D, F), B) = \min(\text{dist}(D, B), \text{dist}(F, B))$$

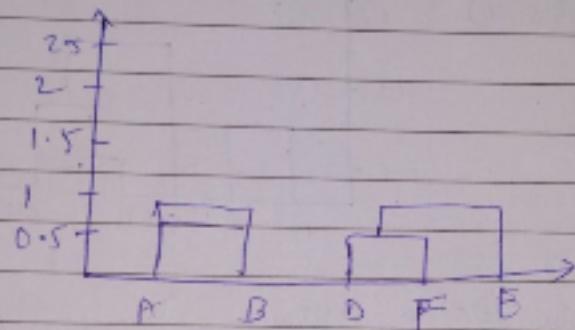
$$= \min(4, 3.5)$$

$$\approx 3.5$$

	A	B	C	(D, F)	E	F
A	0					
B		0				
C	8	7	0			
(D, F)	4.5	3.5	3	0		
E	6	5	2	1	0	
F	4.5	3.5	3.5	0.5		

$$\text{dist} \{(D, F), C\} = \min (\text{dist}(D, C), \text{dist}(F, C)) \\ = \min (3, 3.5) \\ = 3$$

$$\text{dist} \{(D, F), E\} = \min (\text{dist}(D, E), \text{dist}(F, E)) \\ = \min (1, 1.5) \\ = 1$$



Recomputing the distance matrix

$$\text{dist} \{(A, B), C\} \\ = \min (\text{dist}(A, C), \text{dist}(B, C)) \\ = \min (8, 7) \\ = 7.$$

$(A, B) \subset ((D, F), E)$
 $((D, F), B) \begin{bmatrix} 0 & 7 & 0 \\ 7 & 0 & 3.5 \\ 3.5 & 2 & 0 \end{bmatrix}$

$\text{dist}((A,B), ((D,F), E))$

$$= \min [\text{dist}(A,E), \text{dist}(B,E), \text{dist}(A,D), \text{dist}(A,F), \\ \text{dist}(B,D), \text{dist}(BF))]$$

$$= \min (6, 5, 5, 4.5, 4, 3.5)$$

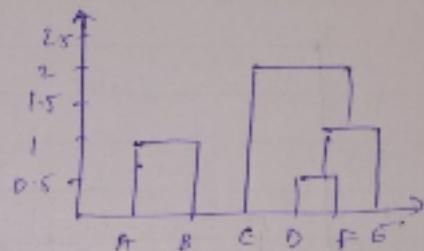
$$= 3.5.$$

$\text{dist}([(A,F), E], E))$

$$= \min [\text{dist}(D,C), \text{dist}(F,C), \text{dist}(E,C)]$$

$$= \min [3, 3.5, 2]$$

$$= 2.$$



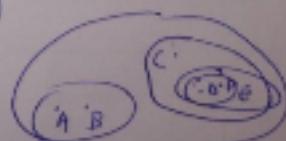
Re Computing the distance matrix $(\underline{(A,B)}, \underline{((D,F), E, C)})$

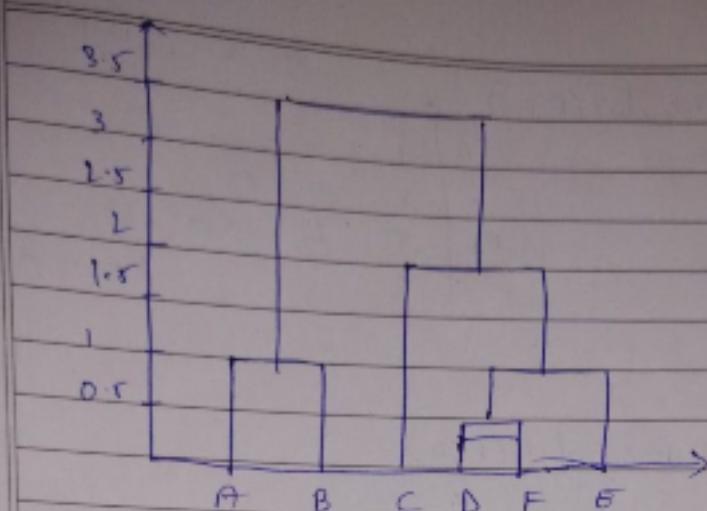
$$\text{dist}((A,B), ((D,F), E, C)) \quad ((D,F), E, C) \begin{bmatrix} 0 \\ 3.5 & 0 \end{bmatrix}$$

$$= \min (\text{dist}(A,B), \text{dist}(A,C), \text{dist}(B,C), \text{dist}(E,C), \\ \text{dist}(A,D), \text{dist}(A,F), \text{dist}(B,D), \text{dist}(B,F))$$

$$= \min (6, 8, 5, 7, 5, 4.5, 4, 3.5)$$

$$= 3.5$$

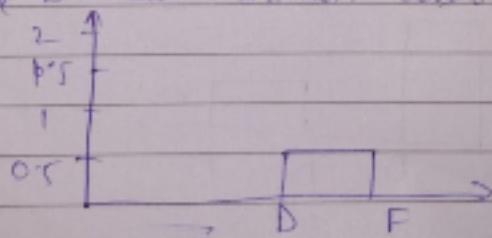




ii) Complete linkage
distance matrix

	A	B	C	D	E	F
A	0					
B	1	0				
C	8	7	0			
D	5	4	3	0		
E	6	5	2	1	0	
F	4.5	3.5	3.5	6.5	18	0

Point F & D has minimum distance 0.5.



Re Computing the distance matrix

$\text{dist}(C, D, F, A)$

$$= \max(\text{dist}(D, A), \text{dist}(F, A))$$

$$= \max(5, 4.5)$$

$$= 5$$

	A	B	C	(D, F)	E
A	0				
B	1	0			
C	8	7	0		
D, F	5	4	3.5	0	
E	6	5	2	1.5	0

$\text{dist}((D, F), B)$

$$= \max(\text{dist}(D, B), \text{dist}(F, B))$$

$$= \max(4, 3.5)$$

$$= 4$$

$$\text{dist}((B, F), C) = \max(\text{dist}(B, C), \text{dist}(F, C))$$

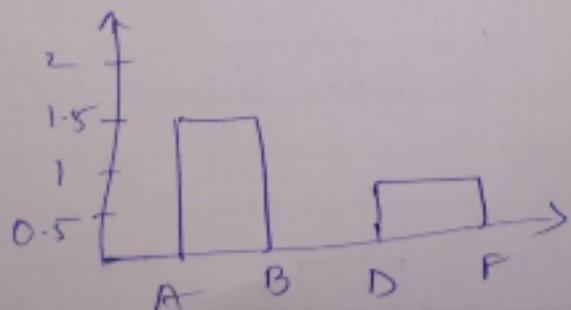
$$= \max(3, 3.5)$$

$$= 3.5$$

$$\text{dist}((D, F), E) = \max(\text{dist}(D, E), \text{dist}(F, E))$$

$$= \max(1, 1.5)$$

$$= 1.5$$



Re computing the distance matrix

$$\text{dist}((A,B), C) = \max(\text{dist}(A,C), \text{dist}(B,C)) \\ = \max(8, 7) \\ = 7$$

$$\begin{array}{c} (A,B) \subset (D,F) \in \\ \begin{array}{c} (A,B) \\ C \\ (D,F) \\ B \end{array} \end{array} \quad \left[\begin{array}{ccc} 0 & & \\ 7 & 0 & \\ 5 & 3.5 & 0 \\ 6 & 2 & 1.5 \end{array} \right]$$

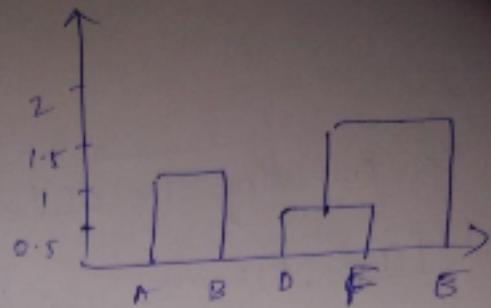
$$\text{dist}((A,B), (A,F)) \\ = \max(\text{dist}(A,D), \text{dist}(A,F), \text{dist}(B,D), \text{dist}(B,F)) \\ = \max(5, 4.5, 4, 3.5) \dots \\ = 5$$

$$\text{dist}((D,F), C) = \max(\text{dist}(D,C), \text{dist}(F,C)) \\ = \max(3, 3.5) \\ = 3.5$$

$$\text{dist}((A,B), E) = \max(\text{dist}(A,E), \text{dist}(B,E)) \\ = \max(6, 5) \\ = 6$$

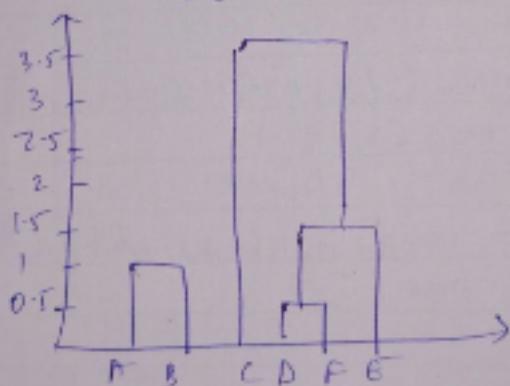
$$\text{dist}((D,F), E) = \max(\text{dist}(D,E), \text{dist}(F,E)) \\ = \max(1, 1.5) \\ = 1.5$$

For Computing the distance
matrix.



$$\begin{aligned}\text{dist}(A, B), \text{dist}(C, D, F), 6) \\ &= \max(\text{dist}(A, E), \text{dist}(B, E), \text{dist}(A, D), \text{dist}(A, F), \\ &\quad \text{dist}(B, D), \text{dist}(B, F)) \\ &= \max(6, 5, 5, 4.5, 9, 3.5) \\ &= 6.\end{aligned}$$

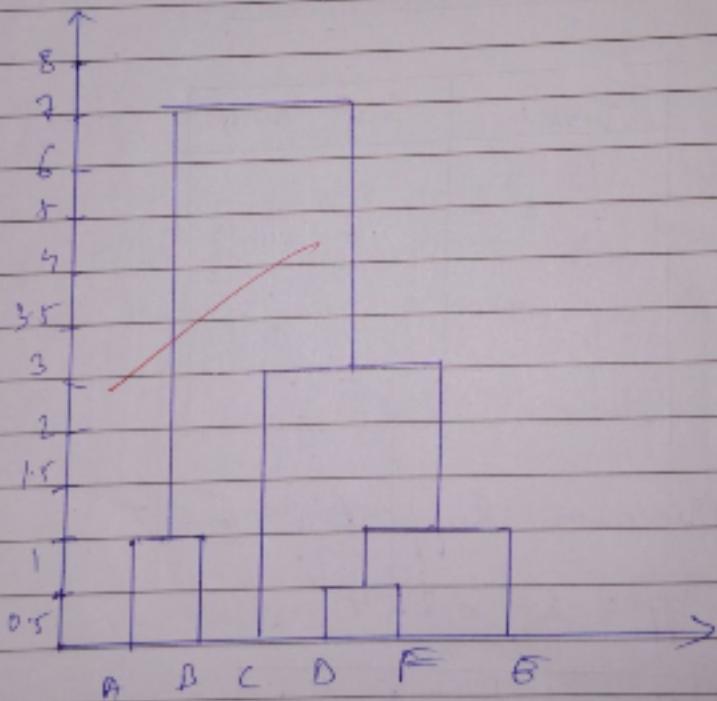
$$\begin{aligned}\text{dist}((C \cup F), E), c) &= \max(\text{dist}(D, c), \text{dist}(F, c), \text{dist}(E, c)) \\ &= \max(3, 3.5, 3) \\ &= 3.5\end{aligned}$$



Precomputing the distance matrix

$$\text{dist}((A,B), (C,D,F), (E,G)) \quad \begin{pmatrix} (A,B) & (C,D,F) & (E,G) \\ (A,B) & 0 & \\ (D,F,E,C) & 8 & 0 \end{pmatrix}.$$

$$\begin{aligned} &= \max (\text{dist}(A,C), \text{dist}(A,G), \\ &\quad \text{dist}(B,E), \text{dist}(B,G) \\ &\quad \text{dist}(A,D), \text{dist}(B,F), \text{dist}(B,D), \text{dist}(B,F)) \\ &= \max (6, 8, 5, 7, 5 + 4, 4, 3) \end{aligned}$$



Q.6. Apply Apriori algorithm of following data to find frequent itemset & strong association rules minimum support = 50% & confidence = 80%

TFID	Items
100	1, 3, 4
200	2, 3, 5
300	1, 2, 3, 5
400	2, 5
500	2, 3, 5
600	2, 5
700	2, 4
800	2, 3, 5

→ Step 1 :-
 C_1

Itemset	Sup-Count
1	2
2	6
3	5
4	2
5	6

Step 2 :- Minimum Support is 50%.

$$\therefore \text{Support} = \frac{50}{100} \times 8 = 4$$

$I_1 =$

Itemset	Sup - Count
2	6
3	5
5	6

Step 3:-

Generate Candidate C₂ from L₁ & find the support of L₁.

C ₂ =	Itemset	Sup-count
	2, 3	4
	2, 5	6
	3, 5	4.

Step 4:-

Compare C₂ generated in step 3 with support count & remove those itemsets which do not satisfy the minimum support count.

L ₂ =	Itemset	Sup-count
	2, 3	4
	2, 5	6
	3, 5	4.

(x)
(x)
(x)

Step 5:-

Generate candidate C₃ from L₂ & find the support.

C ₃ =	Itemset	Sup-count
	2, 3, 5	4

Step 6 :-

Compare G with minimum support count.

Itemset	Sup-count
2,3,5	4.

∴ The database contain the frequent itemsets (2,3,5)

Association rule	Support	Confidence	Confidence P/O.
$2 \wedge 3 \Rightarrow 5$	4	$4/4 = 1$	100.
$3 \wedge 5 \Rightarrow 2$	4	$4/4 = 1$	100
$2 \wedge 5 \Rightarrow 3$	4	$4/6 = 0.667$	66.7.
$5 \Rightarrow 2 \wedge 3$	4	$4/6 = 0.667$	66.7.
$2 \wedge 3 \Rightarrow 5$	4	$4/6 = 0.667$	66.7.
$3 \wedge 2 \Rightarrow 5$	4	$4/6 = 0.667$	66.7.

Given minimum confidence = 80%, so only the first & second rules above are output.

∴ Strong rules are

Rule 1: $2 \wedge 3 \Rightarrow 5$

Rule 2: $3 \wedge 5 \Rightarrow 2$.

These are strong rules because P/O above confidence.