

Regression Analysis :- (gives approx. values)

Regression lines:- ① y on x line \rightarrow This line used to find y when x is given
② x on y line \rightarrow " " " " " " x " y " "

① Regression line y on x defined as

$$y - \bar{y} = b_{yx} (x - \bar{x})$$

② The regression line x on y defined as

$$x - \bar{x} = b_{xy} (y - \bar{y})$$

where,

b_{yx} = Regression coefficient x on y

b_{xy} = " " " " y on x

$$\bar{x} = \frac{\sum x}{n}, \bar{y} = \frac{\sum y}{n}$$

x	y
x_1	y_1
x_2	y_2
:	:
x_n	y_n

Regression Coefficient

$$b_{yx} = \frac{\sum \frac{6y}{6x}}{\sum x} = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\frac{\sum x^2 - (\sum x)^2}{n}}$$

$$\text{where } \sigma_y^2 = \frac{\sum (y - \bar{y})^2}{n}, \sigma_x^2 = \frac{\sum (x - \bar{x})^2}{n}$$

β = coefficient of correlation

$$b_{xy} = \frac{\sum \frac{6x}{6y}}{\sum y} = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\frac{\sum y^2 - (\sum y)^2}{n}}$$

b_{xy} if $u = x - A$ where A = assumed mean of x

$v = y - B$ B = assumed mean of y

$$\text{Then, } b_{xy} = \frac{\sum uv - \frac{\sum u \sum v}{n}}{\frac{\sum v^2 - (\sum v)^2}{n}}, b_{yx} = \frac{\sum uv - \frac{\sum u \sum v}{n}}{\frac{\sum u^2 - (\sum u)^2}{n}}$$

Properties of regression lines

- 1) The regression lines always intersect at the point (\bar{x}, \bar{y})
- 2) If two regression lines are given and θ be the acute angle between them, then show that

$$\tan \theta = \frac{1 - \rho^2}{\rho} \sqrt{\frac{6x \cdot 6y}{6x^2 + 6y^2}}$$

where ρ = coefficient of correlation

- Q) Find the two lines of regression from the following data

Age of husband: 25 22 28 26 35 20 22 40 20 18

Age of wife: 18 15 20 17 22 14 16 21 15 14

- Estimate (i) age of husband when age of wife is 19 and
(ii) age of wife, when age of husband is 30 (iii) correlation coefficient

(Let x = age of husband, y = age of wife and $n = n - 1 = n - 10$,

$v = y - \bar{y} = y - 17$ where 26 and 17 are assumed means
for x and y , and $n = 10$

Here, $\bar{x} = \sum x/n = 256/10 = 25.6$, $\bar{y} = \sum y/n = 172/10 = 17.2$

Now regression coefficient, $b_{xy} = 0.385$

$$b_{xy} = 2.23$$

Given $y = 19$ (age
of wife)

Here the line of regression

$$y \text{ on } x \Rightarrow y - \bar{y} = b_{xy}(x - \bar{x})$$

$$\text{Then, } x = 29.6 \approx 30$$

$$y - 17.2 = 2.23(30 - 25.6)$$

$$\text{or, } y = 0.385x + 7.34 \quad \text{①} \quad \text{② when } x = 30$$

(age of husband)

$$\text{Then } y = 18.89 \approx 19$$

Regression line x on y

$$x - \bar{x} = b_{yx} (y - \bar{y})$$

$$\Rightarrow x = 2.23y - 12.76 \quad \text{②}$$

③ Now the coefficient of

correlation

$$\rho = \sqrt{b_{xy} \cdot b_{yx}}$$

$$= 0.927 \quad \because b_{xy} \text{ & } b_{yx} \text{ are +ve.}$$

Chapter - 6

1. Linear Regression & Correlation:

x	-1	1	2	4	6	7
y	-1	2	3	3	5	8

Sol:-

x	y	x^2	y^2	xy
-1	-1	1	1	1
1	2	1	4	2
2	3	4	9	6
4	3	16	9	12
6	5	36	25	30
7	8	49	64	56

$$\sum x = 19 \quad \sum y = 20 \quad \sum x^2 = 107 \quad \sum y^2 = 112 \quad \sum xy = 107$$

We have;

y-intercept (a) of regression;

$$a = \frac{\sum x \sum xy - \sum y \sum x^2}{(\sum x)^2 - n \sum x^2} = \frac{19 \times 107 - 20 \times 107}{(19)^2 - 6 \times 107} = 0.3808$$

slope (b) of regression;

$$b = \frac{(\sum x)(\sum y) - n \sum xy}{(\sum x)^2 - n \sum x^2} = \frac{19 \times 20 - 6 \times 107}{(19)^2 - 6 \times 107} = 0.9324$$

Now;

$$y = a + b x \quad \text{For } x = -1 \text{ & } y = -1$$

$$\text{or, } -1 = 0.38 + 0.93(-1)$$

$$\therefore y \approx 0.9324 x + 0.3808$$

if x within range - interpolation
 otherwise extrapolation (x & a changes)

$$S_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$= \sqrt{\frac{\sum x^2 - \frac{1}{n} (\sum x)^2}{n-1}}$$

Correlation of: (r)

$$r = \left(\frac{S_x}{S_y} \right) b$$

OR;

$$r = \frac{\sum xy - \frac{1}{n} \sum x \sum y}{(n-1) S_x S_y}$$

Now;

$$S_x = \sqrt{\frac{\sum x^2 - \frac{1}{n} (\sum x)^2}{n-1}} = \sqrt{\frac{107 - \frac{1}{6} (19)^2}{6-1}} = 3.0605$$

$$S_y = \sqrt{\frac{\sum y^2 - \frac{1}{n} (\sum y)^2}{n-1}} = \sqrt{\frac{112 - \frac{1}{6} (20)^2}{6-1}} = 3.0110$$

$$\therefore r = \left(\frac{3.0605}{3.0110} \right) 0.9324 = 0.9477$$

If $y = 4.5$, then

$$\frac{4.5 - 0.3808}{0.9324} = x$$

$$\therefore x = 4.4178 \text{ (interpolation)}$$

2. Multiple Linear Regression:

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$$

where,

α_i = Regression coefficient

x_i = independent variables

y = dependent variables

n = attributes, features

3. Regularization:

Ridge Regression

Lasso Regression

i) Ridge Regression:

$$y = 0.9 + 1.2x_1 + 20x_2 + 39x_3$$

(scale down)

$$\therefore y = 0.9 + 0.7x_1 + 2x_2 + 5x_3$$

ii) Lasso Regression:

$$y = 0.9 + 0x_1 + 0x_2 + 5x_3$$

3. Logistic Regression:

Decision Tree Example:

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	weak	NO
2	Sunny	Hot	High	strong	NO
3	Overcast	Hot	High	weak	Yes
4	Rain	Mild	High	weak	Yes
5	Rain	Cool	Normal	weak	Yes
6	Rain	Cool	Normal	strong	NO
7	Overcast	cool	Normal	strong	Yes
8	Sunny	Mild	High	weak	NO
9	Sunny	Cool	Normal	weak	Yes
10	Rain	Mild	Normal	weak	Yes
11	Sunny	Mild	Normal	strong	Yes
12	Overcast	Mild	High	strong	Yes
13	Overcast	Hot	Normal	weak	Yes
14	Rain	Mild	High	strong	NO

Decision Tree:

1. ID3 (Iterative

$$\text{Entropy}(S) = - \sum p(I) \cdot \log_2 p(I)$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum [p(S|A) \cdot \text{Entropy}(S|A)]$$

where;

S = decision

A = attributes

Here;

9 decisions labelled Yes

5 decisions labelled NO

Total Decision = 14

$$\begin{aligned}\text{Entropy}(\text{Decision}) &= -p(\text{Yes}) \cdot \log_2 p(\text{Yes}) - p(\text{No}) \cdot \log_2 p(\text{No}) \\ &= -(9/14) \cdot \log_2 (9/14) - (5/14) \cdot \log_2 (5/14) \\ &= 0.6428 \cdot 0.6374 + 0.3571 \cdot 1.4854 \\ &= 0.9420\end{aligned}$$

Wind Factor on Decision:

$$\text{Gain}(D, W) = \text{Entropy}(D) - \sum [p(D|W) \cdot \text{Entropy}(D|W)]$$

$$\begin{aligned}\text{Gain}(D, W) &= \text{Entropy}(D) - [p(D|W=\text{Weak}) \cdot \text{Entropy}(D|W=\text{Weak})] \\ &\quad - [p(D|W=\text{strong}) \cdot \text{Entropy}(D|W=\text{strong})]\end{aligned}$$

Weak Wind Factor Decision:

$$\begin{aligned}\text{Entropy (D|W = Weak)} &= -p(\text{No}) * \log_2 p(\text{No}) - p(\text{Yes}) * \log_2 p(\text{Yes}) \\ &= -(2/8) * \log_2 (2/8) - (6/8) * \log_2 (6/8) \\ &= 0.8112\end{aligned}$$

Strong Wind Factor Decision:

$$\begin{aligned}\text{Entropy (D|W = strong)} &= -p(\text{No}) * \log_2 p(\text{No}) - p(\text{Yes}) * \log_2 p(\text{Yes}) \\ &= -(3/6) * \log_2 (3/6) - p(3/6) * \log_2 (3/6) \\ &= 1\end{aligned}$$

Now;

$$\begin{aligned}\text{Gain (D,W)} &= 0.94 - (8/14) * 0.8112 - (6/14) * 1 \\ &= 0.048\end{aligned}$$

Again;

Outlook Factor on Decision:

S: 5, O: 4, R: 5

Now;

$$\begin{aligned}\text{Entropy (D|O = sunny)} &= -(3/5) * \log_2 (3/5) - p(2/5) * \log_2 (2/5) \\ &= 0.9709\end{aligned}$$

$$\begin{aligned}\text{Entropy (D|O = overcast)} &= -(0/4) * \log_2 (0/4) - (4/4) * \log_2 (4/4) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{Entropy (D|O = Rainy)} &= -(2/5) * \log_2 (2/5) - (3/5) * \log_2 (3/5) \\ &= 0.9709\end{aligned}$$

$$\begin{aligned}\text{Gain (D,O)} &= 0.9402 - (5/14) * 0.9709 - (4/14) * 0 - \\ &\quad (5/14) * 0.9709 \\ &= 0.2467\end{aligned}$$

Humidity Factor on Decision:

$$H = 1, N = 1$$

$$\text{Entropy (D|H=High)} = -(4/7) * \log_2(4/7) - (3/7) * \log_2(3/7) \\ = 0.9852$$

$$\text{Entropy (D|H=Normal)} = -(1/7) * \log_2(1/7) - (6/7) * \log_2(6/7) \\ = 0.5916$$

$$\text{Gain (D, H)} = 0.9402 - (7/14) * 0.9852 - (7/14) * 0.5916 \\ = 0.1518$$

Temperature Factor on Decision:

$$H =$$

$$\text{Entropy (D|T=High)} =$$

$$\text{Entropy (D|T=Mild)} =$$

$$\text{Entropy (D|T=cool)} =$$

$$\text{Gain (D, T)} =$$

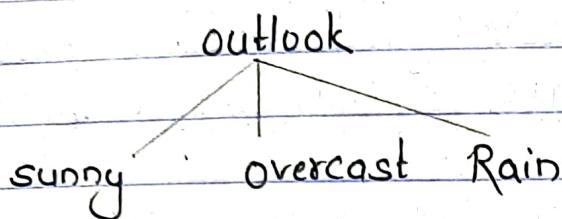
$$= 0.029$$

$$G(D, W) = 0.048$$

$$G(D, H) = 0.151$$

$$G(D, O) = 0.246 \quad \text{Highest}$$

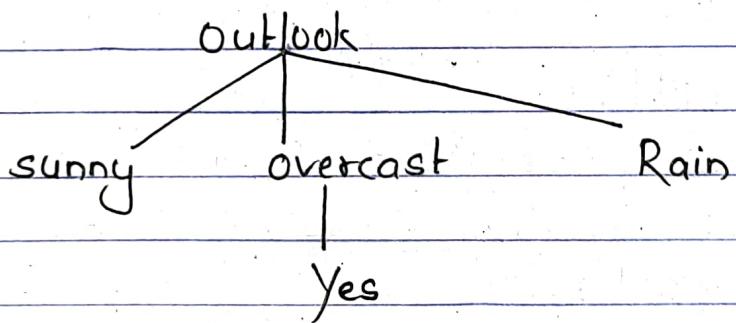
$$G(D, T) = 0.029$$



- outlook factor on decision produces the highest one, so outlook decision will appear in the root node.

Overcast Outlook on decision:

- Decision will always be Yes.



Sunny outlook on decision:

- (outlook = sunny | Temperature) gain = 0.570

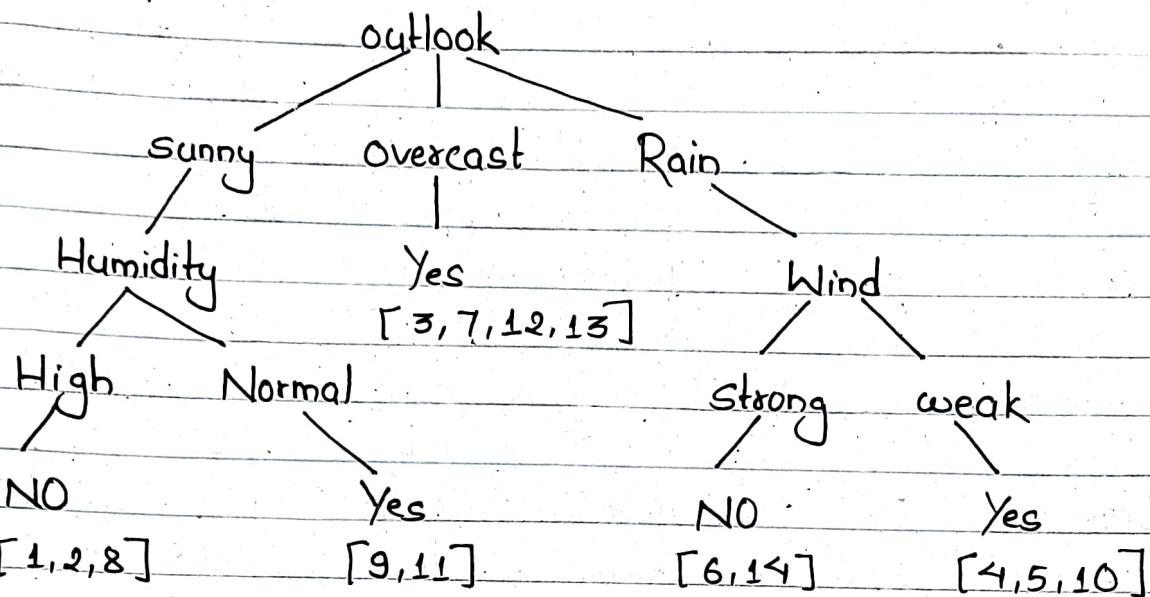
(outlook = sunny | Humidity) gain = 0.970

(outlook = sunny | Wind) gain = 0.019

- Humidity is the decision.

CART, C4.5 & C5.0
(Gini index)

Complete Decision Tree:



Gain (outlook = Rain | Temperature)

Gain (outlook = Rain | Humidity)

Gain (outlook = Rain | Wind)

- Wind has the highest gain.

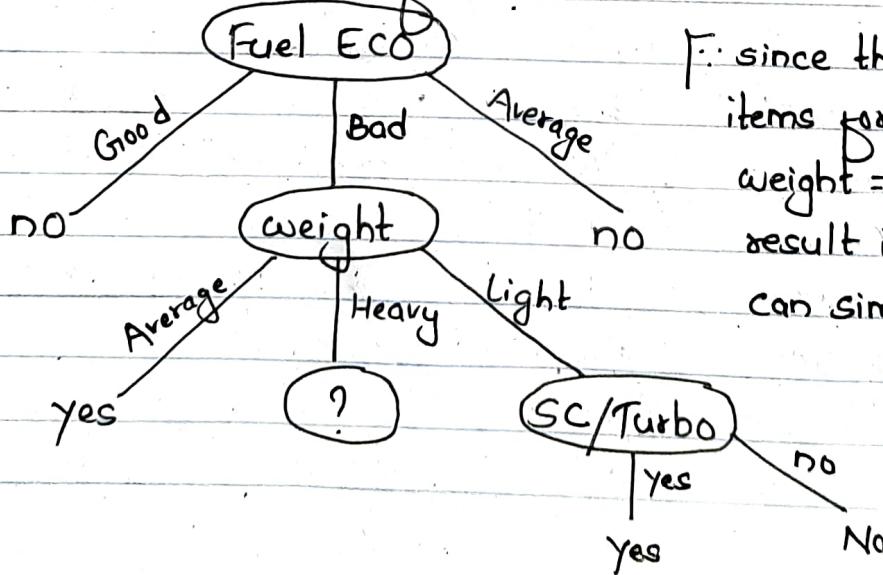
C4.5 Algorithm

- The C4.5 algorithm introduces a number of improvements over the original ID3 algorithm
- The C4.5 algorithm can handle missing data
- If the training records contain unknown attribute values, the C4.5 evaluates the gain for an attribute by considering only the records where the attribute is defined.
- Both categorical and continuous attributes are supported by C4.5
- Values of a continuous variable are sorted and partitioned.
- For the corresponding records of each partition, the gain is calculated, and the partition that maximizes the gain is chosen for the next split.
- The ID3 algorithm may construct a deep & complex tree, which would cause overfitting.
- The C4.5 algorithm addresses the overfitting problem in ID3 by using a bottom-up technique called pruning to simplify the tree by removing the least visited nodes & branches.

2. ID3 Example

Model	Engine	SC/Turbo	Weight	Fuel ECO	Fast
Prius	small	no	average	good	no
Civic	small	no	light	average	no
WRX STI	small	yes	average	bad	Yes
M3	medium	no	heavy	bad	Yes
RSX	large	no	average	bad	Yes
GTI	medium	no	light	bad	no
XJR	large	yes	heavy	bad	no
S500	large	no	heavy	bad	no
911	medium	yes	light	bad	Yes
Corvette	large	no	average	bad	Yes
Insight	small	no	light	good	no
RSX	small	no	average	average	no
IS350	medium	no	heavy	bad	no
MR2	small	yes	average	average	no
E320	medium	no	heavy	bad	no

Is the car fast?



∴ since there are only two items for SC/Turbo where weight = light, and the result is consistent, we can simplify the weight = Light path.]

Bayesian classification:

Bayes Theorem:

$$P(H/X) = \frac{P(X/H) P(H)}{P(X)}$$

→ $P(H/X)$ - posterior probability

Naive Bayesian Classification:

$$P(C_i/X) = \frac{P(X/C_i) P(C_i)}{P(X)}$$

Q.1

RID	age	income	student	credit-rating	buy-computer	Class:
1	Youth	high	no	Fair	no	
2	Youth	high	no	Excellent	no	
3	middle-aged	high	no	Fair	yes	
4	senior	medium	no	Fair	yes	
5	senior	low	yes	Fair	yes	
6	senior	low	yes	Excellent	no	
7	middle-aged	low	yes	Excellent	yes	
8	youth	medium	no	Fair	no	
9	Youth	low	yes	Fair	yes	
10	senior	medium	yes	Fair	yes	
11	Youth	medium	yes	Excellent	yes	
12	middle-aged	medium	no	Excellent	yes	
13	middle-aged	high	yes	Fair	yes	
14	senior	medium	no	Excellent	no	

$X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$
 $\text{buy-computer} = ?$

Soln:- Here, To compute $P(X|c_i)$, for $i=1, 2$

$$P(\text{buy-computer} = \text{yes}) = 9/14 = 0.643$$

$$P(\text{buy-computer} = \text{NO}) = 5/14 = 0.357$$

$$P(\text{age} = \text{youth} | \text{buys-computer} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} | \text{buys-computer} = \text{NO}) = 3/5 = 0.600$$

$$P(\text{income} = \text{medium} | \text{buys-computer} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} | \text{buys-computer} = \text{NO}) = 2/5 = 0.400$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{NO}) = 1/5 = 0.200$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{NO}) = 2/5 = 0.400$$

Now,

$$\begin{aligned} P(X | \text{buys-computer} = \text{yes}) &= P(\text{age} = \text{youth} | \text{buys-computer} = \text{yes}) \\ &\quad \times P(\text{income} = \text{medium} | \text{buys-computer} = \text{yes}) \\ &\quad \times P(\text{student} = \text{yes} | \text{buys-computer} = \text{yes}) \\ &\quad \times P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{yes}) \\ &= 0.222 \times 0.444 \times 0.667 \times 0.667 \\ &= 0.049 \end{aligned}$$

$$\begin{aligned} P(X | \text{buys-computer} = \text{no}) &= 0.600 \times 0.400 \times 0.200 \times 0.400 \\ &= 0.019 \end{aligned}$$

$$\begin{aligned}
 & P(X \mid \text{buys-computer} = \text{yes}) * P(\text{buys-computer} = \text{yes}) \\
 & = 0.044 * 0.693 \\
 & = 0.028
 \end{aligned}$$

$$\begin{aligned}
 & P(X \mid \text{buys-computer} = \text{no}) * P(\text{buys-computer} = \text{no}) \\
 & = 0.019 * 0.357 \\
 & = 0.007
 \end{aligned}$$

Therefore, naive bayesian classifier predicts "buys-computer=yes" for tuple X .

Q.2 (Bayesian Classifier)

Fruit	Yellow	Sweet	Long	Total
Mango	350	450	0	650
Banana	400	300	350	400
others	50	100	50	150
Total	800	850	400	1200

$$\text{Fruit} = \{ \text{Yellow, sweet, Long} \}$$

Soln:- Here;

$$\begin{aligned}
 P(\text{Fruit} = \text{Mango}) &= 650/1200 = 0.541 \\
 P(X \mid \text{Mango}) &= P(Y|M) * P(S|M) * P(L|M)
 \end{aligned}$$

$$\begin{aligned}
 P(\text{Yellow} \mid \text{Mango}) &= \frac{P(\text{Mango} \mid \text{Yellow}) \cdot P(\text{Yellow})}{P(\text{Mango})} \\
 &= \frac{350/800 * 800/1200}{650/1200} \\
 &= \frac{350}{650} \\
 &= 0.53
 \end{aligned}$$

i) Mango

$$\begin{aligned}
 P(X \mid \text{Mango}) &= P(Y/M) \cdot P(S/M) \cdot P(L/M) \\
 &= 0.53 * \frac{450/850 * 850/1200}{650/1200} * \frac{900/1200 * 0}{650/1200} \\
 &= 0.53 * 0.692 * 0 \\
 &= 0
 \end{aligned}$$

$$\therefore P(S/M) = 0.692$$

$$\therefore P(L/M) = 0$$

$$\therefore P(X \mid \text{Mango}) = 0$$

ii) Banana

$$P(\text{Banana}) = 900/1200 = 0.75$$

$$P(Y|B) = \frac{P(B|Y) * P(Y)}{P(B)} = \frac{\frac{500}{800} * \frac{800}{1200}}{\frac{500}{1200}} = 1$$

$$P(S|B) = \frac{P(B|S) * P(S)}{P(B)} = 0.75$$

$$P(L|B) = 0.875, \text{ then}$$

$$P(X|B) = P(Y|B) * P(S|B) * P(L|B)$$

$$\therefore P(X|B) = 0.65$$

iii) others:

$$P(Y|0) = 0.33$$

$$P(S|0) = 0.66$$

$$P(L|0) = 0.33 \text{ then,}$$

$$P(X|0) = P(Y|0) * P(S|0) * P(L|0)$$

$$\therefore P(X|0) = 0.072$$

∴ $P(X|B)$ has higher probability value.

∴ Fruit = {Yellow, sweet, long} = Banana //

Association Rule:

- Support
- Confidence

$$\text{Support} = \frac{(X \cup Y) \cdot \text{count}}{n}$$

$$\text{Support} = \frac{n \cdot x}{n}$$

$$\text{Confidence} = \frac{(X \cup Y) \cdot \text{count}}{X \cdot \text{count}}$$

$$\text{Confidence} = \frac{\text{support}(X \cup Y)}{\text{support}(X)}$$

$$\text{Support} = \frac{2}{9} \times 100\% = 22.22\%$$

Apriori Algorithm:

1. For the following given transaction dataset generate rules using Apriori Algorithm. Consider the values as minimum support = $\frac{2}{9}$ and minimum confidence = 70%.

Transaction ID	List of Items
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

Soln:- Here,

C_1 (Candidates)

Item-set	Support Count (frequency)
{A}	6
{B}	7
{C}	6
{D}	2
{E}	2

L_1 (Support ≥ 2)

Item Set	Support count
$\{A\}$	6
$\{B\}$	7
$\{C\}$	6
$\{D\}$	2
$\{E\}$	2

C_2 (Candidates pair)

Item Set	Support count
$\{A, B\}$	4
$\{A, C\}$	4
$\{A, D\}$	1
$\{A, E\}$	2
$\{B, C\}$	4
$\{B, D\}$	2
$\{B, E\}$	2
$\{C, D\}$	0
$\{C, E\}$	1
$\{D, E\}$	0

L_2 (Support ≥ 2)

Item Set	Support Count
$\{A, B\}$	4
$\{A, C\}$	4
$\{A, E\}$	2
$\{B, C\}$	4
$\{B, D\}$	2
$\{B, E\}$	2

C₃

Item set	Support Count
{A, B, C}	2
{A, B, E}	2
{A, C, E}	1
{B, C, D}	0
{B, C, E}	1
{B, D, E}	0
{A, B, D}	1

L₃ (support ≥ 2)

Item set	Support count
{A, B, C}	2
{A, B, E}	2

C₄

Item set	Support Count
{A, B, C, E}	1
{A, B, C, D}	
{B, C, E, D}	

L₄ (support ≥ 2) = \emptyset

Since $L_4 = \emptyset$, we take selected combination i.e {A, B, C} and {A, B, E} i.e we consider rule for L_3 .

For Combination {A, B, C} :

x	y	Confidence
{A, B} \rightarrow {C}		$A \cup B \cup C / A \cup B = 2/4 = 50\%$
{A, C} \rightarrow {B}		$2/4 = 50\%$
{B, C} \rightarrow {A}		$2/4 = 50\%$
{A} \rightarrow {B, C}		$2/6 = 33.33\%$
{B} \rightarrow {A, C}		$2/7 = 28.57\%$
{C} \rightarrow {A, B}		$2/6 = 33.33\%$

i.e $\frac{\text{support}(x \cup y)}{\text{support}(x)}$

Since all the rules do are less than 70% of confidence, we don't accept these rules.

For the confidence of combination $\{A, B, E\}$

$\{A, B\} \rightarrow \{E\}$	Confidence
$\{A, E\} \rightarrow \{B\}$	$A \cup B \cup E / A \cup B = 2/4 = 50\%$
$\{B, E\} \rightarrow \{A\}$	$2/2 = 100\%$
$\{A\} \rightarrow \{B, E\}$	$2/6 = 33.33\%$
$\{B\} \rightarrow \{A, E\}$	$2/7 = 28.57\%$
$\{E\} \rightarrow \{A, B\}$	$2/2 = 100\%$

Final Association Rules; from above calculations, we get following rules: (as their confidence is greater than 70%)

$\{A, E\} \rightarrow \{B\}$
 $\{B, E\} \rightarrow \{A\}$
 $\{E\} \rightarrow \{A, B\}$ //

Apriori property \rightarrow All non-empty subsets of a frequent item sets must also be frequent.

2. Generate Association Rules using Apriori Algorithm.

Minimum Support = 3

Minimum Confidence = 70%

Transaction ID	List of Items
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, U, C, K, Y}
T500	{C, O, O, K, I, E}

Soln:- Here,

C ₁ (Candidates)		L ₁ (Support ≥ 3)	
Item Set	Support Count	Item Set	Support Count
{M}	3	{M}	3
{O}	3	{O}	3
{N}	2	{K}	5
{K}	5	{E}	4
{E}	4	{Y}	3
{Y}	3		
{D}	1		
{A}	1		
{U}	1		
{C}	2		
{O _r }	1		
{I}	1		

C_2 (Candidate Pairs)

Item set	Support Count
$\{M, O\}$	1
$\{M, K\}$	3
$\{M, E\}$	2
$\{M, Y\}$	2
$\{O, K\}$	3
$\{O, E\}$	3
$\{O, Y\}$	2
$\{K, E\}$	4
$\{K, Y\}$	3
$\{E, Y\}$	2

L_2 (Support ≥ 3)

Item set	Support Count
$\{M, K\}$	3
$\{O, K\}$	3
$\{O, E\}$	3
$\{K, E\}$	4
$\{K, Y\}$	3

C_3 (Candidates)

Item set	Support Count
$\{M, O, K\}$	1
$\{K, O, E\}$	3
$\{K, E, Y\}$	2
$\{M, K, E\}$	2
$\{M, K, Y\}$	2
$\{O, K, Y\}$	2

L_3 (Support ≥ 3)

Item set	Support count
$\{K, O, E\}$	3

We select combination $\{K, O, E\}$ for association rule.

For Confidence of Combination $\{K, O, E\}$

	Confidence
$\{K, O\} \rightarrow \{E\}$	$KOUE / KOO = 3/3 = 100\%$
$\{K, E\} \rightarrow \{O\}$	$3/4 = 75\%$
$\{O, E\} \rightarrow \{K\}$	$3/3 = 100\%$
$\{K\} \rightarrow \{O, E\}$	$3/5 = 60\%$
$\{O\} \rightarrow \{K, E\}$	$3/3 = 100\%$
$\{E\} \rightarrow \{K, O\}$	$3/4 = 75\%$

Hence the required final association rule from above calculations are as follows:

$$\begin{aligned}\{K, O\} &\rightarrow \{E\} = 100\% \\ \{K, E\} &\rightarrow \{O\} = 75\% \\ \{O, E\} &\rightarrow \{K\} = 100\% \\ \{O\} &\rightarrow \{K, E\} = 100\% \\ \{E\} &\rightarrow \{K, O\} = 75\%\end{aligned} //$$

3. Support = 50%, Confidence = 75%

TID	Items Purchased
1	Bread, Cheese, Egg, Juice
2	Bread, Cheese, Juice
3	Bread, Milk, Yogurt
4	Bread, Juice, Milk
5	Cheese, Juice, Milk

Soln:- Here, C_1 (candidates)

Item Set	Support Count	Support
Bread	4	$4/5 = 80\%$
Cheese	3	$3/5 = 60\%$
Egg	1	$1/5 = 20\%$
Juice	4	$4/5 = 80\%$
Yogurt	1	$1/5 = 20\%$
Milk	3	$3/5 = 60\%$

L_1 (support $\geq 50\%$)

Item set	Support Count
Bread	80%
Cheese	60%
Juice	80%
Milk	60%

C_2 (Candidate pairs)

Item set	Support Count	Support
{B, C}	2	$2/5 = 40\%$
{B, J}	3	$3/5 = 60\%$
{B, M}	2	$2/5 = 40\%$
{C, J}	3	$3/5 = 60\%$
{C, M}	1	$1/5 = 20\%$
{J, M}	2	$2/5 = 40\%$

L_2 (support $\geq 50\%$)

Item set	Support
{B, J}	60%
{C, J}	60%

C_3

$$l_3 (\text{support} \Rightarrow 50\%) = \emptyset$$

Itemset	Support Count	Support
$\{B, C, J\}$	1	$1/5 = 20\%$

Since $l_3 = \emptyset$, we take selected combination i.e $\{B, J\}$ and $\{C, J\}$ i.e rule. L_2 .

For confidence of combination $\{B, J\}$

	confidence
$\{B\} \rightarrow \{J\}$	$B \cup J / B = 3/4 = 75\%$

For confidence of combination $\{C, J\}$

	confidence
$\{C\} \rightarrow \{J\}$	$C \cup J / C = 3/3 = 100\%$

Hence the required final association rule from the above calculation are as follows:

$$\{B\} \rightarrow \{J\} = 75\%$$

$$\{C\} \rightarrow \{J\} = 100\%$$

||

4. Minimum Support = 30%

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

Soln:- Here,

C₁ (candidates)

Item set	Support Count	Support
{A}	5	5/8 = 62.5%
{B}	6	6/8 = 75%
{C}	3	3/8 = 37.5%
{D}	6	6/8 = 75%
{E}	4	4/8 = 50%

L₁ (support = 30%)

Item set	Support
{A}	62.5%
{B}	75%
{C}	37.5%
{D}	75%
{E}	50%

C₂ (candidate Pairs)

Item set	Support count	support
{A, B}	4	4/8 = 50%
{A, C}	2	2/8 = 25%
{A, D}	4	4/8 = 50%
{A, E}	4	4/8 = 50%
{B, C}	3	3/8 = 37.5%

C_2 (Candidate Pairs)

Item Set	Support count	support
$\{B, D\}$	4	$4/8 = 50\%$
$\{B, E\}$	3	$3/8 = 37.5\%$
$\{C, D\}$	1	$1/8 = 12.5\%$
$\{C, E\}$	2	$2/8 = 25\%$
$\{D, E\}$	3	$3/8 = 37.5\%$

L_2 (support $\Rightarrow 30\%$)

Item set	support
$\{A, B\}$	50%
$\{A, D\}$	50%
$\{A, E\}$	50%
$\{B, C\}$	37.5%
$\{B, D\}$	50%
$\{B, E\}$	37.5%
$\{D, E\}$	37.5%

C_3 (Candidates)

Item Set	Support count	support
$\{A, B, D\}$	3	$3/8 = 37.5\%$
$\{A, D, E\}$	3	$3/8 = 37.5\%$
$\{A, B, C\}$	2	$2/8 = 25\%$
$\{B, C, D\}$	1	$1/8 = 12.5\%$
$\{D, B, E\}$	2	$2/8 = 25\%$
$\{B, C, E\}$	2	$2/8 = 25\%$
$\{A, B, E\}$	3	$3/8 = 37.5\%$

L_3 (support $\Rightarrow 30\%$)

Item set	Support
$\{A, B, D\}$	37.5%
$\{A, D, E\}$	37.5%
$\{A, B, E\}$	37.5%

C_4 (candidates)

Item Set	Support count	support
$\{A, B, D, E\}$	2	$2/8 = 25\%$

$$L_4 = \emptyset$$

Since $L_4 = \emptyset$, we take selected combination i.e $\{A, B, D\}$, $\{A, D, E\}$ and $\{A, B, E\}$ i.e we consider rule for L_3 .

FP Growth:

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

Minimum Support = 3

Soln:- Here,

C_1 (candidates)

Item Set	Support Count
{M}	3
{O}	3
{N}	2
{K}	5
{E}	4
{Y}	3
{D}	1
{A}	1
{U}	1
{C}	2
{I}	1

L_1 (Support ≥ 3)

Item set	Support count
{M}	3
{O}	3
{K}	5
{Y}	3
{E}	4

Items	Ordered items with high frequency	Priority	
{K}	5	1	
{E}	4	2	
{M}	3	3	
{O}	3	4	
{Y}	3	5	

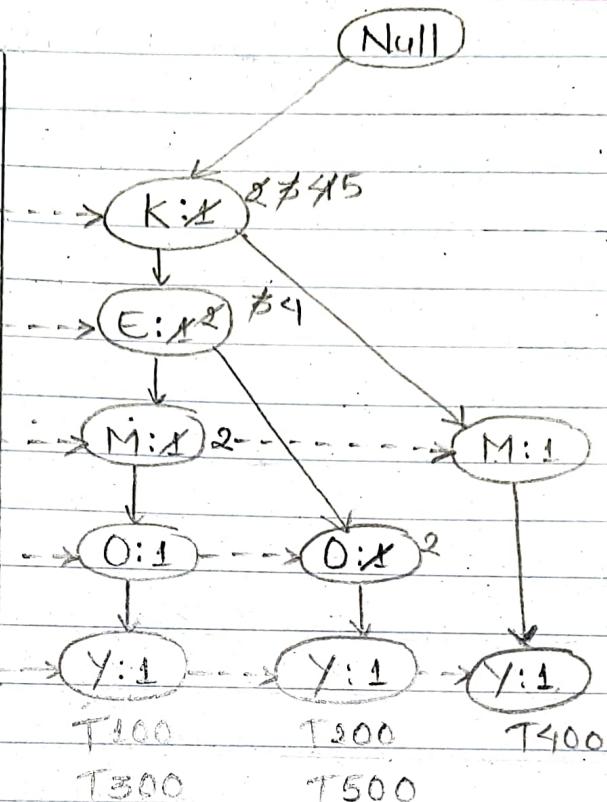
TID	Items	ordered items with priority
T100	{M, O, N, K, E, Y}	{K, E, M, O, Y}
T200	{D, O, N, K, E, Y}	{K, E, O, Y}
T300	{M, A, K, E}	{K, E, M}
T400	{M, U, C, K, Y}	{K, M, Y}
T500	{C, O, K, I, E}	{K, E, O}

Generating FP Tree:

Freq. item	Support count
{K}	5
{E}	4
{M}	3
{O}	3
{Y}	3

FP Tree:

freq. Item	Support count
$\{K\}$	5
$\{E\}$	4
$\{M\}$	3
$\{O\}$	3
$\{Y\}$	3



Mining FP Tree:

Items (low priority)	Conditional Pattern Base	Conditional FP Tree	Frequent Patterns Generated
$\{Y\}$	$\{\{K, E, M, O:1\}, \{K, E, O:1\}, \{K, M:1\}\}$	$\langle K:3 \rangle$	$\{K, Y:3\}$
$\{O\}$	$\{\{K, E, M:1\}, \{K, E:2\}\}$	$\langle K:3, E:3 \rangle$	$\{K, O:3\}, \{E, O:3\}$
$\{M\}$	$\{\{K, E:2\}, \{K:1\}\}$	$\langle K:3 \rangle$	$\{K, M:3\}$
$\{E\}$	$\{K:4\}$	$\langle K:4 \rangle$	$\{K, E:4\}$
$\{K\}$	-		

$$2. \text{ Minimum Support} = 30\% = \frac{30}{100} \times 8 = 2.4 \approx 3$$

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

Soln:- Here;

C_1 (candidates)

Item set	Support count
{A}	5
{B}	6
{C}	3
{D}	6
{E}	4

I_1 (support ≥ 3)

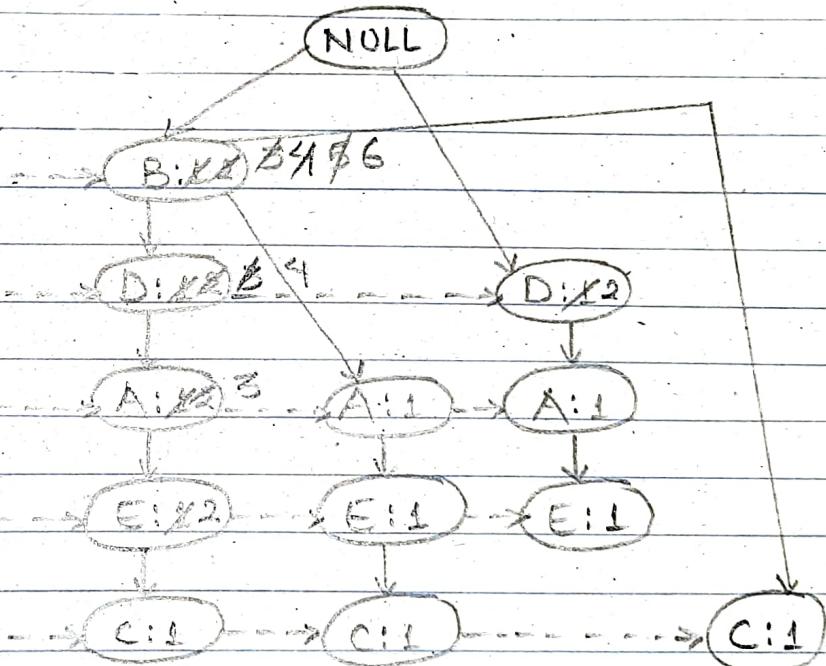
Item set	Support
{A}	5
{B}	6
{C}	3
{D}	6
{E}	4

Items	ordered items with high freq.	Priority
B	6	1
D	6	2
A	5	3
E	4	4
C	3	5

TID	Items	ordered items with priority
1	{E, A, D, B}	{B, D, A, E}
2	{D, A, C, E, B}	{B, D, A, E, C}
3	{C, A, B, E}	{B, A, E, C}
4	{B, A, D}	{B, D, A}
5	{D}	{D}
6	{D, B}	{B, D}
7	{A, D, E}	{D, A, E}
8	{B, C}	{B, C}

Generating FP Tree:

Frequent Item	Support Count
B	6
D	6
A	5
E	4
C	3



Mining FP Tree:

Item (Low priority ↓)	Conditional Pattern Base	Conditional FP Tree	frequent Pattern Generated
$\{C\}$	$\{\{B, D, A, E: 1\}, \{B, A, E: 1\}, \{B: 1\}\}$	$\langle B: 3 \rangle$	$\{B, C: 3\}$
$\{E\}$	$\{\{B, D, A: 2\}, \{B, A: 1\}, \{D, A: 1\}\}$	$\langle B: 3, D: 3, A: 1 \rangle$	
$\{A\}$	$\{\{B, D: 3\}, \{B: 1\}, \{D: 1\}\}$	$\langle B: 4, D: 1 \rangle$	$\{B, A: 4\}, \{D, A: 3\}$
$\{D\}$			$\{B, D, A: 3\}$
$\{B\}$		$\langle B: 4 \rangle$	$\{B, D: 4\}$

K-mean: (One point)

1. Using K-mean clustering algorithm, divide the given dataset into 2 clusters ($K=2$).

Data	$ 4-2 =2$	$ 12-2 =10$	$ m-d $
	$m_1 = 9$	$m_2 = 12$	
2	2	10	m_1
3	1	9	m_1
4	0	8	m_1
10	6	2	m_2
11	7	1	m_2
12	8	0	m_2
20	16	8	m_2
25	21	13	m_2
30	26	18	m_2

Soln:- Let $m_1 = 9$ & $m_2 = 12$.

Then;

Clusters formed are;

$$K_1 = \{2, 3, 4\} \quad \text{and} \quad K_2 = \{10, 11, 12, 20, 25, 30\}$$

$$m_1(\text{new}) = \frac{2+3+4}{3} = 3 \quad m_2 = \frac{10+11+12+20+25+30}{6} = 18$$

Iteration 2:

$$K_1 = \{2, 3, 4, 10\}$$

$$m_1(\text{new}) = 4.75 \approx 5$$

$$K_2 = \{11, 12, 20, 25, 30\}$$

$$m_2 = 19.6 \approx 20$$

Iteration 3:

data	$m_1 = 5$	$m_2 = 20$	
2	3	18	m_1
3	2	17	m_1
4	1	16	m_1
10	5	10	m_1
11	6	9	m_1
12	7	8	m_1
20	15	0	m_2
25	20	5	m_2
30	25	10	m_2

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$m_1(\text{new}) = 7$$

$$K_2 = \{20, 25, 30\}$$

$$m_2(\text{new}) = 25$$

Iteration 4:

Data	$m_1 = 7$	$m_2 = 25$	
2	5	23	m_1
3	4	22	m_1
4	3	21	m_1
10	3	15	m_1
11	4	14	m_1
12	5	13	m_1
20	13	5	m_2
25	18	0	m_2
30	23	5	m_2

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

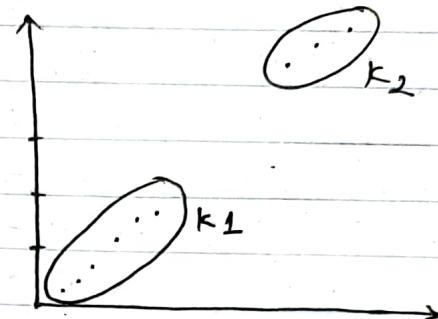
$$K_2 = \{20, 25, 30\}$$

Since the clusters k_1 & k_2 formed in iteration 4 are same as of iteration 3, so we stop here.

∴ Final Clusters are;

$$k_1 = \{2, 3, 4, 10, 11, 12\}$$

$$k_2 = \{20, 25, 30\}$$



2. Use k-mean algorithm;

Data = {15, 15, 16, 19, 19, 20, 20, 21, 22, 28, 35, 40, 41, 42, 43, 44, 60, 61, 65}

$k = 2$ and centroid ($c_1 = 16$, $c_2 = 22$)

Soln:- Here;

$$K = 2$$

$$c_1 = m_1 = 16$$

$$c_2 = m_2 = 22$$

Iteration 1:

Data	$ m_1 - d $	$ m_2 - d $	
15	1	7	m_1
15	1	7	m_1
16	0	6	m_1
19	3	3	m_1
19	3	3	m_1
20	1	2	m_2
20	1	2	m_2
21	5	1	m_2
22	6	0	m_2
28	12	6	m_2
35	19	13	m_1
40	24	18	m_2
41	25	19	m_2
42	26	20	m_2
43	27	21	m_2
44	28	22	m_2
60	44	38	m_2
61	45	39	m_2
65	49	43	m_2

$$K_1 = \{15, 15, 16, 19, 19\}$$

$$K_2 = \{20, 20, 21, 22, 28, 35, 40, 41, 42, 43, 44, 60, 61, 65\}$$

$$m_1(\text{new}) =$$

$$m_2(\text{new}) =$$

Minkowski Distance:

$$d(i, j) = q \sqrt{|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{in} - x_{jn}|^q}$$

If $q=1$, d is Manhattan distance,

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$$

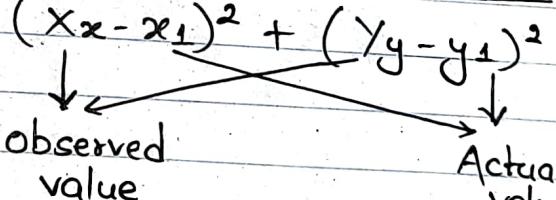
If $q=2$, d is Euclidean distance,

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{in} - x_{jn}|^2}$$

Euclidean Distance:

$$(x, y) \& C(x, y)$$

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


observed value Actual value

Two Point K-mean:

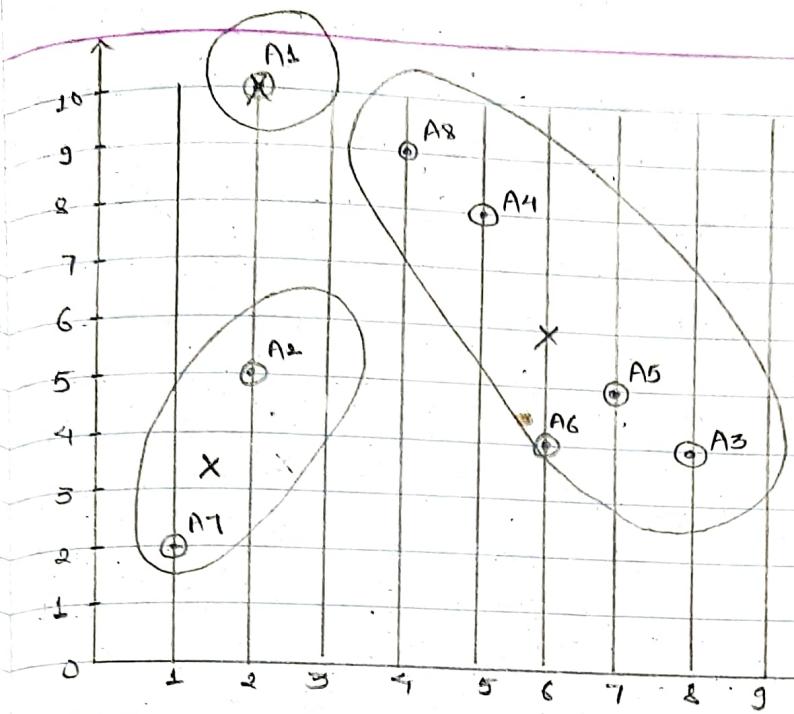
1. Cluster the following 8 points (with (x, y) representing locations) into 3 clusters $A_1(2, 10)$, $A_2(2, 5)$, $A_3(8, 4)$, $A_4(5, 8)$, $A_5(7, 5)$, $A_6(6, 4)$, $A_7(1, 2)$ and $A_8(4, 9)$. Initial cluster centers are: $A_1(2, 10)$, $A_4(5, 8)$ and $A_7(1, 2)$. The distance function between two points $a = (x_1, y_1)$ and $b = (x_2, y_2)$ is defined as:
- $$d(a, b) = |x_2 - x_1| + |y_2 - y_1|$$
- Use k-means algorithm to find the three cluster centers.

Soln:- Here;

Iteration 1:

Point	Distance mean 1 (2, 10)	Distance mean 2 (5, 8)	Distance mean 3 (1, 2)	Cluster
$A_1(2, 10)$	0	5	9	$m_1 / 1$
$A_2(2, 5)$	5	8	4	$m_3 / 3$
$A_3(8, 4)$	12	7	9	m_2
$A_4(5, 8)$	5	0	10	m_2
$A_5(7, 5)$	10	5	9	m_2
$A_6(6, 4)$	10	5	7	m_2
$A_7(1, 2)$	9	10	0	m_3
$A_8(4, 9)$	3	2	10	m_2

Cluster 1	Cluster 2	Cluster 3
$A_1(2, 10)$	$A_3(8, 4)$	$A_2(2, 5)$
$A_4(5, 8)$		$A_7(1, 2)$
$A_5(7, 5)$		
$A_6(6, 4)$		
$A_8(4, 9)$		
$new(2, 10)$	$new(6, 6)$	$new(1.5, 3.5)$

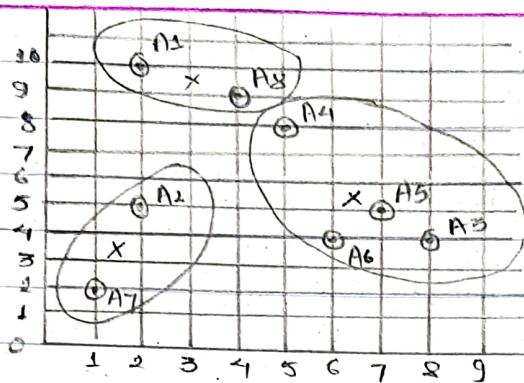


Iteration 2:

Point	mean ₁ (2, 10)	mean ₂ (6, 6)	mean ₃ (1.5, 3.5)	Cluster
A ₁ (2, 10)	0	8	7	m ₁
A ₂ (2, 5)	5	5	2	m ₃
A ₃ (8, 4)	12	4	7	m ₂
A ₄ (5, 8)	5	3	8	m ₂
A ₅ (7, 5)	10	2	7	m ₂
A ₆ (6, 4)	10	2	5	m ₂
A ₇ (1, 2)	9	9	2	m ₃
A ₈ (4, 9)	3	5	8	m ₁

Cluster 1	cluster 2	cluster 3
A ₁ (2, 10)	A ₃ (8, 4)	A ₂ (2, 5)
A ₈ (4, 9)	A ₄ (5, 8)	A ₇ (1, 2)

New mean:
 cluster 1 = (3, 9.5)
 cluster 2 = (6.5, 5.25)
 cluster 3 = (1.5, 3.5)



Iteration 3:

Point	mean ₁ (3.9, 5)	mean ₂ (7.5, 3.5)	mean ₃ (2.5, 4.75)	Cluster
A ₁ (2, 10)	1.5	9.25	7	m ₁
A ₂ (2, 5)	5.5	4.75	2	m ₃
A ₃ (8, 4)	10.5	2.75	7	m ₂
A ₄ (5, 8)	3.5	4.25	8	m ₁
A ₅ (7, 5)	8.5	0.75	7	m ₂
A ₆ (6, 4)	8.5	1.75	5	m ₂
A ₇ (1, 2)	9.5	8.75	2	m ₃
A ₈ (4, 9)	1.5	6.25	8	m ₁

A₉

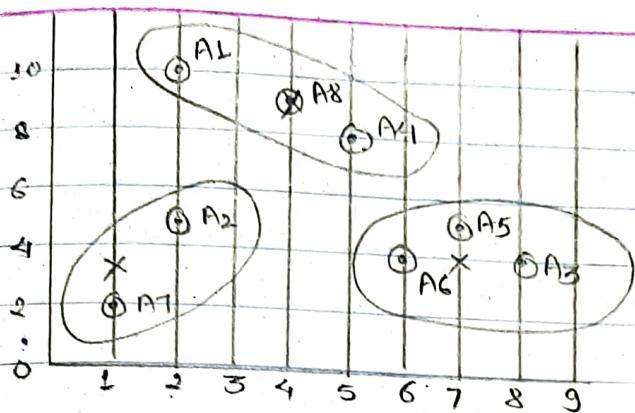
cluster 1	cluster 2	cluster 3
A ₁ (2, 10)	A ₃ (8, 4)	A ₂ (2, 5)
A ₄ (5, 8)	A ₅ (7, 5)	A ₆ (6, 4)
A ₈ (4, 9)	A ₆ (6, 4)	A ₇ (1, 2)

New mean:

$$\text{mean } 1 = (3.67, 9) \approx (4, 9)$$

$$\text{mean } 2 = (7, 4.33) \approx (7, 4)$$

$$\text{mean } 3 = (2.5, 4.75)$$



Iteration 4:

Point	mean ₁ (4,9)	mean ₂ (7,4)	mean ₃ (1.5,3.5)	Cluster
A ₁ (2,10)	3	11	7	m ₁
A ₂ (2,5)	6	6	2	m ₃
A ₃ (8,4)	9	1	7	m ₂
A ₄ (5,8)	2	6	8	m ₁
A ₅ (7,5)	7	1	7	m ₂
A ₆ (6,4)	7	1	5	m ₂
A ₇ (1,2)	10	8	2	m ₃
A ₈ (4,9)	0	8	8	m ₁

Cluster 1	cluster 2	cluster 3
A ₁ (2,10)	A ₃ (8,4)	A ₂ (2,5)
A ₄ (5,8)	A ₅ (7,5)	A ₇ (1,2)
A ₈ (4,9)	A ₆ (6,4)	

$$\text{mean } 1 = (4,9)$$

$$\text{mean } 2 = (7,4)$$

$$\text{mean } 3 = (1.5,3.5)$$

Since the clusters formed in iteration 3 & 4 are same, we stop here.

2. Using K-mean clustering algorithm, divide the given dataset into 2 clusters. Use Euclidean distance.

ID	X	Y
1	1	1
2	1.5	2
3	3	4
4	5	7
5	3.5	5
6	4.5	5
7	3.5	4.5

$$E.C.D = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Soln:- Here, Iteration 1:

ID	X	Y	ECD ₁ (3,4)	ECD ₂ (4.5,5)	Cluster
1	1	1	3.60	5.31	ECD ₁
2	1.5	2	2.5	4.24	ECD ₁
3	3	4	0	1.80	ECD ₁
4	5	7	3.6	2.06	ECD ₂
5	3.5	5	1.11	1	ECD ₁
6	4.5	5	1.80	0	ECD ₂
7	3.5	4.5	0.70	1.11	ECD ₁

Cluster 1	Cluster 2
1 (1,1)	4 (5,7)
2 (1.5,2)	6 (4.5,5)
3 (3,4)	
5 (3.5,5)	
7 (3.5,4.5)	

$$\text{new Mean1} = (2.5, 3.3)$$

$$\text{new Mean2} = (4.5, 6)$$

Iteration 2:

ID	X	Y	ECD ₁ (2.5, 3.3)	ECD ₂ (4.5, 5)	Cluster
1	1	1	2.74	5.31	m_2 ECD ₁
2	1.5	2	1.64	4.24	ECD ₁
3	3	4	0.86	1.80	ECD ₁
4	5	7	4.46	2.06	ECD ₂
5	3.5	5	1.97	1	ECD ₂
6	4.5	5	2.62	0	ECD ₂
7	3.5	4.5	1.56	1.11	ECD ₂

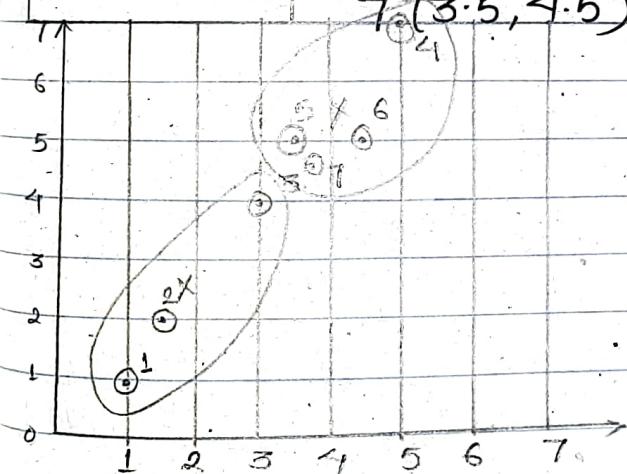
Cluster 1

1 (1, 1)
2 (1.5, 2)
3 (3, 4)

Cluster 2

4 (5, 7)
5 (3.5, 5)
6 (4.5, 5)
7 (3.5, 4.5)

new Mean₁ = (1.8, 2.3)
new Mean₂ = (4.25, 5.37)



Iteration 3:

ID	X	Y	ECD ₁ (1.8, 2.3)	ECD ₂ (4.12, 5.37)	Cluster
1	1	1	1.5	5.4	ECD ₁
2	1.5	2	0.4	4.3	ECD ₁
3	3	4	2.1	1.8	ECD ₂
4	5	7	5.6	1.9	ECD ₂
5	3.5	5	3.2	0.7	ECD ₂
6	4.5	5	3.8	0.5	ECD ₂
7	3.5	4.5	2.8	1.1	ECD ₂

Cluster 1

1 (1, 1)

2 (1.5, 2)

cluster 2

3 (3, 4)

4 (5, 7)

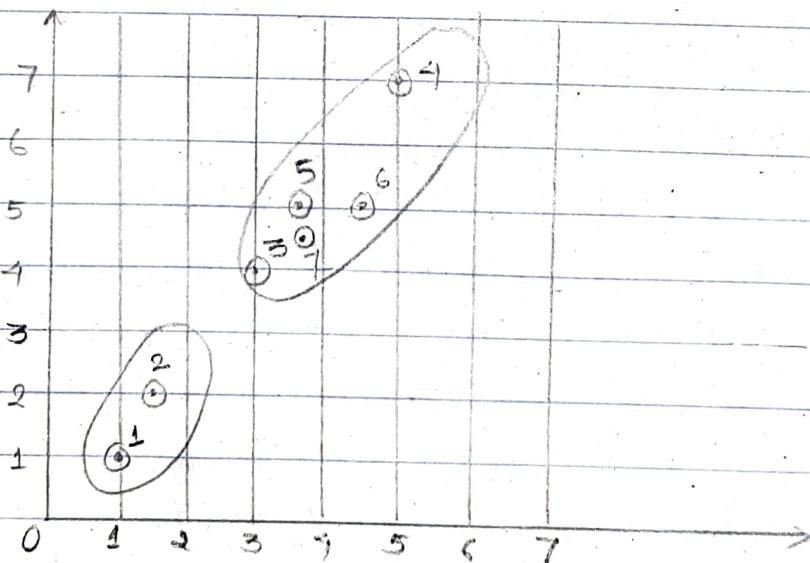
5 (3.5, 5)

6 (4.5, 5)

7 (3.5, 4.5)

New mean₁ = (1.25, 1.5)

New mean₂ = (3.9, 5.1)



Iteration 4:

ID	X	Y	ECD ₁ (1.25, 1.5)	ECD ₂ (3.9, 5.1)	Cluster
1	1	1			
2	1.5	2	0.5	5	ECD ₁
3	3	4	0.5	3.9	ECD ₁
4	5	7	3.1	1.4	ECD ₂
5	3.5	5	6.7	2.2	ECD ₂
6	4.5	5	4.2	0.9	ECD ₂
7	3.5	4.5	4.8	0.6	ECD ₂
			3.8	0.7	ECD ₂

Cluster 1	cluster 2
1 (1, 1)	3 (3, 4)
2 (1.5, 2)	4 (5, 7)
	5 (3.5, 5)
	6 (4.5, 5)
	7 (3.5, 4.5)

Since the clusters formed in iteration 3 & 4 are same, we stop here.

3. Using k-means clustering algorithm, divide the given dataset into 2 clusters. Use Euclidean distance.

Height	Weight
185	72
170	56
168	60
179	68
182	72
188	77
180	71
180	80 70
183	84
180	88
180	67
177	76

K-Medoid:

i	x	y
x_1	2	6
x_2	3	4
x_3	3	8
x_4	4	7
x_5	6	2
x_6	6	4
x_7	7	3
x_8	7	4
x_9	8	5
x_{10}	7	6

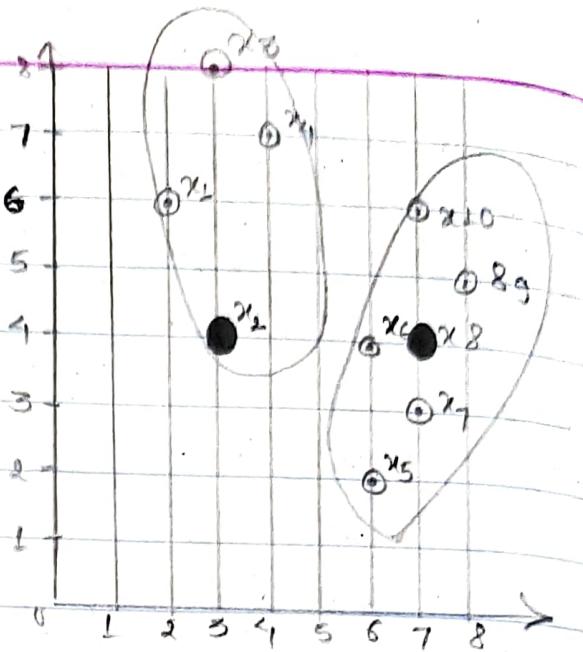
Take $k=2$ and use Manhattan distance.

Soln:- Here, Step 1:

Iteration 1:

i	x	y	Distance for med1 (3,4)	Distance for med2 (7,4)	clusters
x_1	2	6	3	7	med1
x_2	3	4	0	4	med1
x_3	3	8	4	8	med1
x_4	4	7	4	6	med1
x_5	6	2	5	3	med2
x_6	6	4	3	1	med2
x_7	7	3	5	1	med2
x_8	7	4	4	0	med2
x_9	8	5	6	2	med2
x_{10}	7	6	5	2	med1

Cluster 1	Cluster 2
$x_1(2,6)$	$x_5(6,2)$
$x_2(3,4)$	$x_6(6,4)$
$x_3(3,8)$	$x_7(7,3)$
$x_4(4,7)$	$x_8(7,4)$
	$x_9(8,5)$
	$x_{10}(7,6)$



Step 2:

Calculate the total cost, $T_{\text{cost}}(x, c) = \sum_{i=1}^d |x_i - c_i|$

$$\begin{aligned}
 \text{Current Total cost} &= \text{cost } \{(2,6), (3,4)\} + \text{cost } \{(3,9), (3,8)\} + \\
 &\quad \text{cost } \{(3,4), (4,7)\} + \text{cost } \{(7,4), (6,2)\} + \\
 &\quad \text{cost } \{(7,9), (6,4)\} + \text{cost } \{(7,4), (7,3)\} + \\
 &\quad \text{cost } \{(7,4), (8,5)\} + \text{cost } \{(7,4), (7,6)\} \\
 &= (3+4+4) + (3+1+1+2+2) \\
 &= 20
 \end{aligned}$$

Step 3:

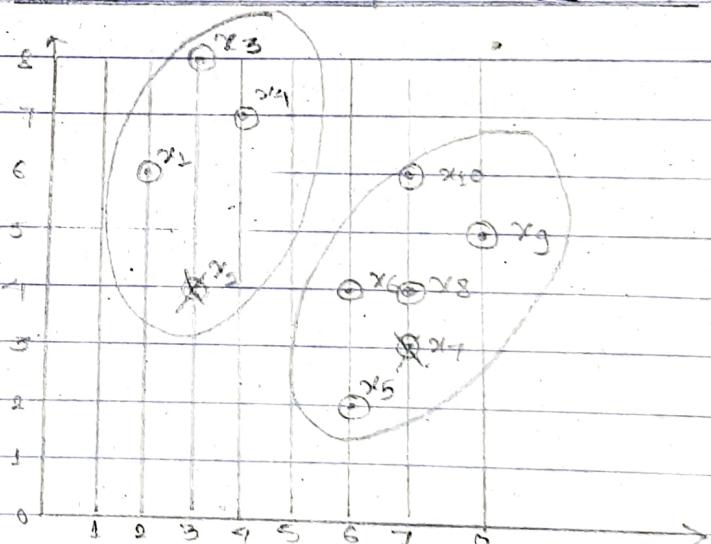
Select one of non-medoid O' and let O' be $(7,3)$.

Step 4:

Iteration 2:

i	x	y	Distance for $me_1(3,4)$	Distance for $me_2(7,3)$	Clusters
x_1	2	6	3	8	me_1
x_2	3	4	0	5	me_1
x_3	3	8	4	9	me_1
x_4	4	7	4	7	me_1
x_5	6	2	5	2	me_2
x_6	6	4	3	2	me_2
x_7	7	3	5	0	me_2
x_8	7	4	4	1	me_2
x_9	8	5	6	3	me_2
x_{10}	7	6	6	3	me_2

Cluster 1	Cluster 2
$x_1(2,6)$	$x_5(6,2)$
$x_2(3,4)$	$x_6(6,4)$
$x_3(3,8)$	$x_7(7,3)$
$x_4(4,7)$	$x_8(7,4)$
$x_9(8,5)$	
$x_{10}(7,6)$	



Step 5:

$$\begin{aligned}
 \text{Current Total Cost} &= (3+4+4) + \text{cost}\{(7,3), (6,2)\} + \\
 &\quad \text{cost}\{(7,3), (6,4)\} + \text{cost}\{(7,3), (7,4)\} + \\
 &\quad \text{cost}\{(7,3), (8,5)\} + \text{cost}\{(7,3), (7,6)\} \\
 &= 11 + (2+2+1+3+3) \\
 &= 22
 \end{aligned}$$

Step 6:

Now, cost of swapping medoid m_2 to $0'$;

$$s = \text{total cost } (m_1, 0') - \text{total cost } (m_1, m_2)$$
$$= 22 - 20$$
$$= 2 > 0$$

So moving $0'$ would be a bad idea as $s > 0$. Hence previous choice was good and the final clusters with medoids $(3, 4)$ and $(7, 4)$ were taken.

2. Take $k=2$ and use Manhattan distance for following dataset:
 $\{(8,7), (3,7), (4,9), (9,6), (8,5), (5,8), (7,3), (8,4), (7,5), (4,5)\}$.

Use K-Medoid Method / algorithm.

Example of K-Medoids:

1. PAM (Partitioning Around Medoids) (For small datasets)
 2. CLARA (clustering LArge Applications)
 3. CLARANS (clustering LArge Applications based upon Randomized Search)
- } For large dataset

K-Mean Vs. K-Medoid:

K-Medoid

1. The complexity of each iteration is, $O(K(n-k)^2)$.

For large values of n & k , such computation becomes very costly.

2. K-Medoid method is more robust than K-Means in the presence of noise and outliers.

3. K-Medoid is more costly.

4. Like K-Means, K-Medoids require the user to specify k .

K-Mean

1. The complexity of each iteration is, $O(K(n-k)^2)$.

For large values of n & k , such computation becomes very costly.

2. K-Mean is less robust than K-Medoid.

3. K-Mean is less costly.

4. K-Mean requires the user to specify k .

K-Medoids Vs K-Means:

- The complexity of each iteration is $O(k(n-k)^2)$.
- For large values of n and k , such computations becomes very costly.

Advantages:

- K-Medoid method is more robust than K-Means in the presence of noise and outliers.

Disadvantages:

- K-Medoid is more costly than the K-Mean method.
- Like K-Mean, k-Medoid requires the user to specify k .
- It does not scale well for large data sets.

Hierachial Method:

1. Agglomerative
2. Divisive

1. Perform agglomerative algorithm on the following data and plot a dendrogram using single link approach. The given data indicates the distances between elements.

Elements	E	A	C	B	D	
E	0					Pair (E,A)
A	1	0				$\Rightarrow (1)$
C	2	2	0			
B	2	5	1	0		
D	3	3	6	3	0	

Soln:- Here;

	(E,A)	C	B	D	
(E,A)	0				Pair (B,C)
C	2	0			$\Rightarrow (1)$
B	2	1	0		
D	3	6	3	0	

	(E,A)	(B,C)	D	
(E,A)	0			Pair ((E,A), (B,C))
(B,C)	2	0		$\Rightarrow (2)$
D	3	3	0	

	(E, A, B, C)	D
(E, A, B, C)	0	
D	3	0

Plotting into Dendrogram;

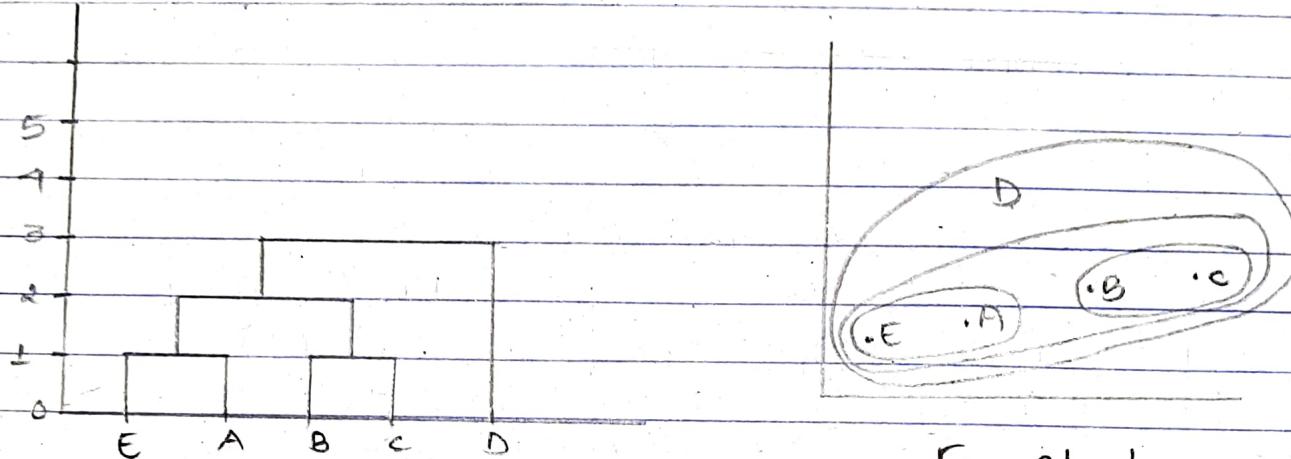


Fig: Cluster

Fig: Dendogram

Item	A	B	C	D	E
A	0				
B	9	0			
C	3	7	0		
D	6	5	9	0	
E	11	10	2	8	0

Pair (C, E)
 $\Rightarrow (2)$

Soln:- Here;

Items	(C,E)	A	B	D
(C,E)	0			
A	3	0		
B	7	9	0	
D	8	6	5	0

Pair (C,E,A)
 $\Rightarrow (3)$

Items	(C,E,A)	B	D
(C,E,A)	0		
B	9	7	0
D	6	5	0

Pair (B,D)
 $\Rightarrow (5)$

Items	(A,C,E)	(B,D)
(A,C,E)	0	
(B,D)	6	0

$$\begin{aligned}
 & \min((A,8), (C,8), (E,8), \\
 & (A,10), (C,9), (E,9)) \\
 & = \min(8, 7, 0, 6, 9, 8) \\
 & = \min(6)
 \end{aligned}$$

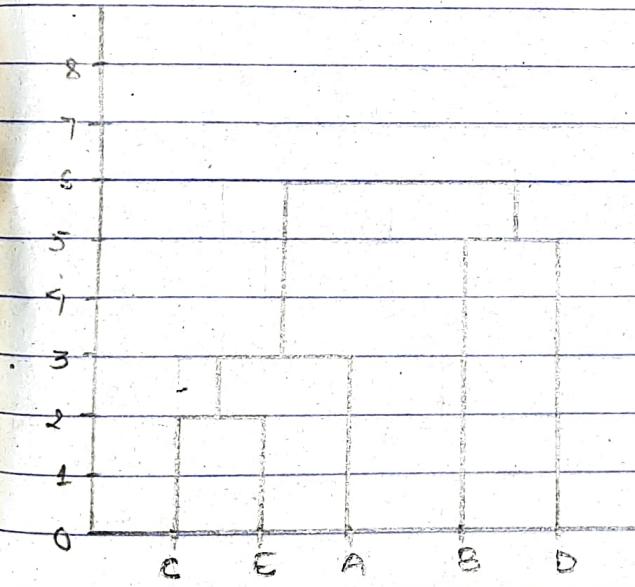


Fig: Dendrogram

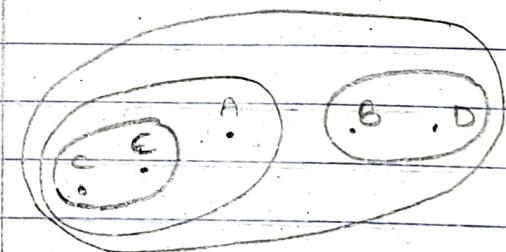


Fig: Clusters

Elements	A	B	C	D	E	F
A	0					
B	0.71	0				
C	5.66	9.95	0			
D	3.61	2.92	2.24	0		
E	4.24	3.54	1.41	1.00	0	
F	3.20	2.50	2.50	0.50	1.12	0

Soln:- Here;

Since pair (D, F) has minimum value 0.50, we select (D, F) as pair.

\downarrow Pair (D, F)
(0.50)

Elements	(D, F)	A	B	C	E
(D, F)	0				
A	3.20	0			
B	2.50	0.71	0		
C	2.24	5.66	9.95	0	
E	1.00	4.24	3.54	1.41	0

For $(D, F) \& (D, F)$;

$$= \min((D, D), (D, F), (F, D), (F, F))$$

$$= \min(0, 0.50, 0)$$

$$= 0$$

For $(D, F) \& A$;

$$= \min((D, A), (F, A))$$

$$= (3.61, 3.20)$$

$$= 3.20$$

The pair (B, A) has the minimum value of 0.71.

Elements	(B, A)	(D, F)	C	E
(B, A)	0			
(D, F)	2.50	0		
C	4.95	2.24		
E	3.54	1.00	1.41	0

\Downarrow pair (D, E, F)
 (1.00)

Elements	(A, B)	C	(D, E, F)
(A, B)	0		
C	4.95	0	
(D, E, F)	2.50	1.41	0

\Downarrow pair (C, D, E, F)
 (1.41)

Elements	(A, B)	(C, D, E, F)
(A, B)	0	
(C, D, E, F)	2.50	0

Plotting into Dendogram;

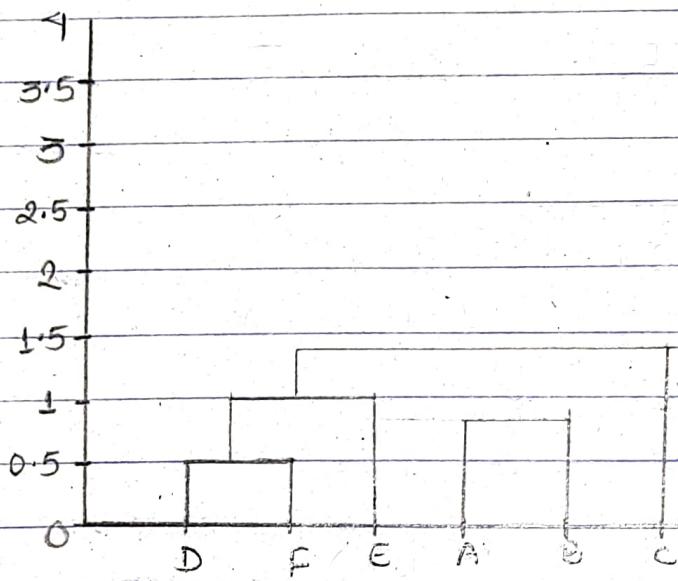


Fig: Dendogram.

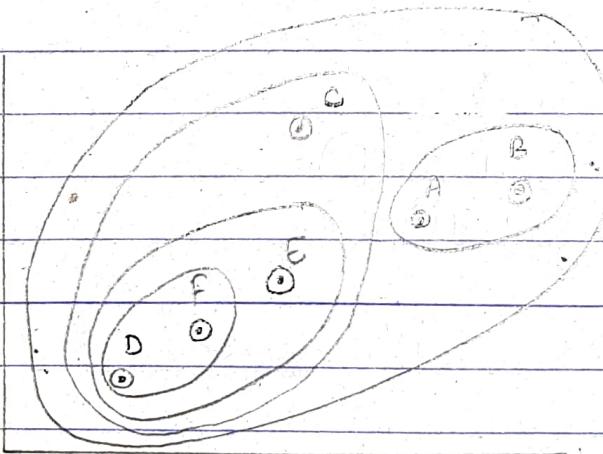


Fig: Clusters.