- Initialization diagram:
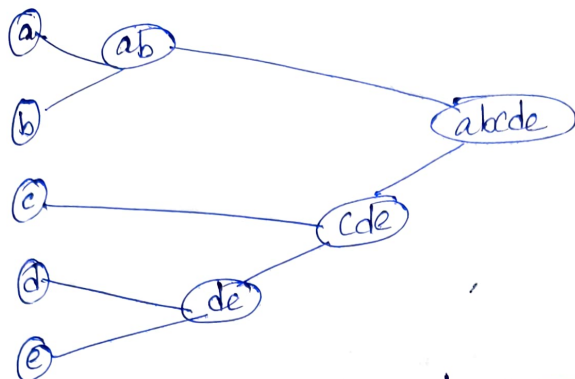
step 0    step 1    step 2    step 3    step 4

→ agglomerative (AGNES)



← 
step4    step3    step2    step1    step 0

divisive (DIANA)

hierarchical clustering

| → Agglomerative clustering | Divisive clustering |
|---|---|
| • bottom up approach | • top - down approach |
| • starts by placing each obj. in its own clusters and then merges these atomic clusters, until all the objects are in a single cluster or until certain termination conditions are satisfied. | • start at the top with all documents in one cluster. <br> • The cluster is split using a flat clustering algorithm. <br> • This procedure is applied recursively until each document is in its own singleton cluster. |

14/3/25

=> Decision tree:

It is an ~~un~~supervised learning technique that gives the input data to be built upon some ML algorithms and predicting the output, based upon the decisions. eg: Yes/No
male or female, class A/B/c.

| eg: Tid | Refund | Marital status | Taxable income | cheat |
|---|---|---|---|---|
| | | categ | conti | class |
| 1 | Yes | Single | 125 K | No |
| 2 | No | Married | 100 K | No |
| 3 | No | Single | 70 K | No |
| 4 | Yes | Married | 120 K | No |
| 5 | No | Divorced | 95 K | Yes |
| 6 | No | Married | 60 K | No |
| 7 | Yes | Divorced | 220 K | No |
| 8 | No | Single | 85 K | Yes |
| 9 | No | Married | 75 K | No |
| 10 | No | Single | 90 K | Yes |

# Splitting Attributes

Refund

No

- splitting attribute based upon the root element. After finding the class, again split the attributes based on subnodes.
- After splitting the process based on subnodes, after finding the class category, going to take good decisions.

eg:

| Name | Rank | Years | Tenured |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant prof | 7 | yes |

→ start from root node –

Refund

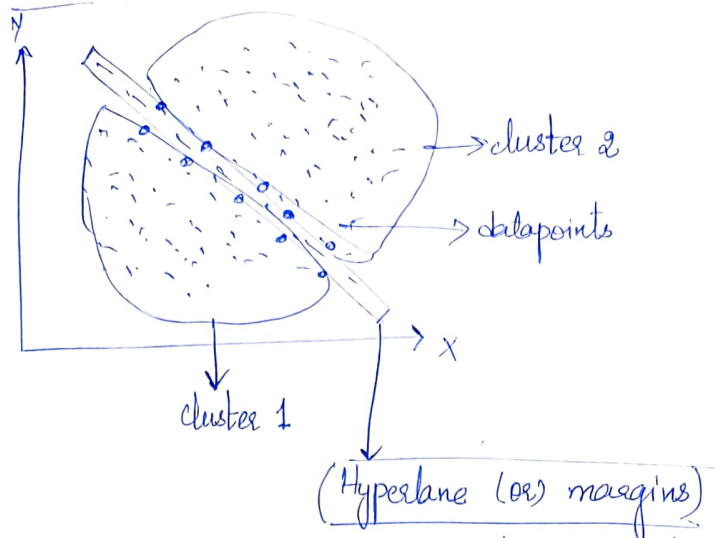| Refund | marital status | Taxable income | cheat |
|--------|----------------|----------------|-------|
| | | | ? |
| No | Married | 80k | |

Refund

No    Mastut

⇒ classification:

- It is a supervised learning technique that predicts the value based on categorial data. It also uses, label the data to predict the output.
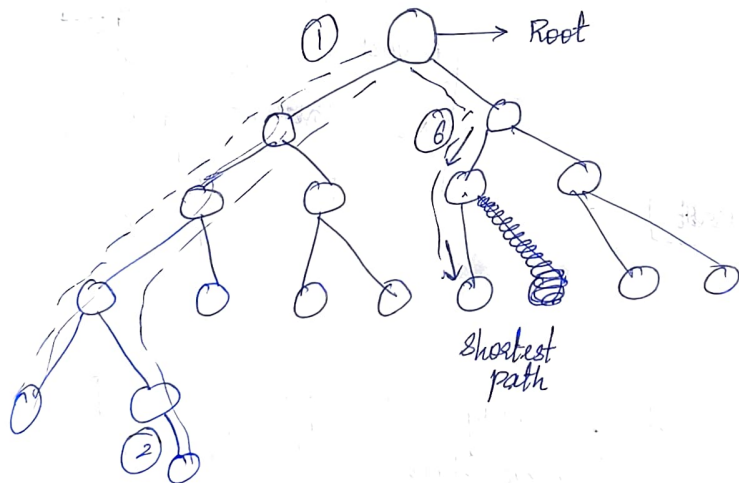
- It classifies data based on testing data.

eg: Naives Bayes, Bayesian network.

Bayesian network → classification algorithus

training data → classifier (model)

if rank = 'professor'
OR years > 6
THEN Tenured = yes

# Support vector machine :



cluster 2
datapoints
cluster 1

(Hyperlane (or) margins)

The dataset is divided into 2 parts through a separate single straight line

In the hyperlane the datasets of different groups like lie on the margin to divide the partitions as segments of the whole dataset.

# Random Forest :



Root

Shortest path

largest path → ②
shortest path → ⑥ → optimized path
(choosing the - optimize solutions.)
shortest path

⇒ Decision tree Algorithm Problem:

Also known as ID3 Algorithm – Iterative Dichotomiser Algorithm.

S:

| | Size | color | shape | Class |
|---|---|---|---|---|
| 1 | Small | Yellow | Round | A |
| 2 | Big | Yellow | Round | A |
| 3 | Big | Red | Round | A |
| 4 | Small | Red | Round | A |
| 5 | Small | Black | Round | B |
| 6 | Big | Black | Cube | B |
| 7 | Big | Yellow | cube | B |
| 8 | Big | Black | round | B |
| 9 | Small | Yellow | Cube | B |
| 10 | | | | |

↑ find no. of classes

Given:

$$M = 2$$

no. of category or class

class A = 4
class B = 5
Total no. of classes = 9

measure of how much uncertainity in the info
entrophy – in the info
Info gain – how much info is the answer to a specific question provided.

Information Gain Entrophy.

Formula for entrophy :-

$$Info(D) = -\sum_{i=1}^{M} P_i \log_2 P_i \rightarrow category$$

↓ Decision

class A  class B

$P_i = \dfrac{no.\ of\ category}{Training\ data}$

class A = 4/9
class B = 5/9

$$Info(D) = \left[-\sum_{i=1} \frac{4}{9} \log \frac{4}{9}\right] + \left[\sum_{i=1}^{2} \frac{5}{9} \log \frac{5}{9}\right]$$

$$= \left[(-0.44)(1.18)\right] + \left[(-0.55) \cdot (-0.86)\right]$$

$$= 0.519 + 0.473 = 0.992$$

$$\left(Info(D) = 0.992\right)//$$

$\text{Info (shape)} =$

$M = 2$

no. of Round (A) = 4

Round (B) = 2

$$\text{Info(D)} = -\sum_{i=1}^{M} P_i \log P_i$$

$$\text{Info(shape)} = \underset{\substack{\uparrow \\ \text{total training data}}}{\frac{\overset{\text{total round}}{\overbrace{6}}}{9}} \left[ \left( -\frac{4}{6} \log_2 \frac{4}{6} \right) + \left( -\frac{2}{6} \log_2 \frac{2}{6} \right) \right]$$

class A      round    class B

$$+ \frac{\overset{3}{\overbrace{3}}}{\underset{\text{total training data}}{9}} \left[ \left( -\frac{0}{3} \log_2 \frac{0}{3} \right) + \left( -\frac{3}{3} \log_2 \frac{3}{3} \right) \right]$$

total cube             cube

$$= \frac{2}{3} \left[ \left( -\frac{2}{3} \log_2 \frac{2}{3} \right) + \left( -\frac{1}{3} \log_2 \frac{1}{3} \right) \right]$$

$$+ \frac{1}{3} \left[ 0 + \left( -1 \times \log_2 (1) \right) \right]$$

$$= \frac{2}{3} \left[ \left( (-0.66)(-0.59) \right) + \left( (-0.33)(-1.59) \right) \right] + 0.33 \left[ -1 \times 0 \right]$$

$$= 0.66 \left[ 0.39 + 0.52 \right] = \underline{0.6006}$$

$\text{Info (color)} = ?$

Total no. of color = 3

~~Yellow (A)~~ = 2

no. of class = 3

$$\text{info(D)} = -\sum_{i=1}^{M} P_i \log P_i$$

class A (yellow)        class B (yellow)

$$\text{Info (color)} = \frac{4}{9} \left[ \left( -\frac{2}{4} \log_2 \frac{2}{4} \right) + \left( -\frac{2}{4} \log_2 \frac{2}{4} \right) \right]$$

$$+ \frac{3}{9} \left[ \left( -\frac{0}{3} \log_2 \frac{0}{3} \right) + \left( \frac{3}{3} \log_2 \frac{-3}{3} \right) \right] + \frac{2}{9} \left[ \left( -\frac{2}{2} \log_2 \frac{2}{2} \right) + \left( -\frac{0}{2} \log_2 \frac{0}{2} \right) \right]$$

Entropy highest value.

→ root → highest value in entropy



Take 2 attribute

[final decision tree]

color, shape, round

class A

class B

⇒ Bayesion classification:

⇒ Decision tree algorithm:

Example 2:

| Owns home | Married | Gender | Employee | class |
|-----------|---------|--------|----------|-------|
| | | | Y | B |
| Y | Y | M | Y | A |
| Y | N | F | Y | C |
| N | Y | F | N | B |
| Y | N | M | Y | C |
| Y | Y | F | Y | C |
| N | N | F | Y | A |
| N | N | M | Y | B |
| N | N | M | N | A |
| Y | Y | F | Y | C |
| N | Y | F | | |
| Y | Y | F | | |

no. of classes = 3 (A, B, C)
(M)

22/3/25

⇒ **Bayesion classification:**

| outlook | Temperature | Humidity | Windly | class |
|---|---|---|---|---|
| sunny | Hot | high | False | N |
| sunny | Hot | high | True | N |
| overcast | Hot | high | False | P |
| rain | Mild | high | False | P |
| rain | Cool | Normal | False | P |
| rain | Cool | Normal | True | N |
| overcast | Cool | normal | True | P |
| sunny | mild | high | False | N |
| sunny | cool | normal | False | P |
| Rain | Mild | Normal | False | P |
| sunny | Mild | Normal | True | P |
| Overcast | ~~hot~~ mild | high | True | P |
| overcast | ~~mild~~ hot | Normal | False | P |
| rain | mild | high | True | N |

**Rule:**

**Step 1** $P(x) =$ outlook $=$ ① Rain, temp $=$ ② hot, Humidity $=$ ③ high, windy $=$ ④ false

No. of class $= 2$

Total no. of training data $= 14$

$$P(N) = \frac{5}{14} = 0.357$$

$$P(P) = \frac{9}{14} = 0.642$$

$P(\text{outlook}) = p$ rain $/ N.) = \frac{2}{5} = 0.4$ → Total no. of N

$P(\text{outlook}) =$ rain $/ p) = \frac{3}{9} = 0.33$ → Total no. of P.

$p(temp = hot/N) = \frac{2}{5} = 0.4$

$p(temp = hot/P) = \frac{2}{9} = 0.22$

$p(humidity = high/N) = \frac{4}{5} = 0.8$

$p(humidity = high/P) = \frac{3}{9} = 0.33$

$p(windy = false/N) = \frac{2}{5} = 0.4$

$p(windy = false/P) = \frac{6}{9} = 0.66$

step 2:

Formula $\Rightarrow$ $P(H/x) = \frac{P(x/H) \cdot P(H)}{P(x)} = P(x_1 c_1) \times p(x_2 c_2)$

$\cdots P(x_n c_n)$

$P(x/N) = P(N/x)\, p(outlook = rain/N) \times p(temp = hot/N)$
$\times p(humid = high/N) \times p(windy = false/N)$

$= 0.4 \times 0.4 \times 0.8 \times 0.4 = 0.0512$

$P(x/P) = p(outlook = rain/P) \times p(temp = hot/P) \times p(humid = high/P)$
$\times p(windy = false/N)$

$= 0.3^{a} \times 0.2 \times 0.3 \times 0.6 = \underline{0.0108}$

at last step:

$p(x/N) = 0.0512$

$P(x/P) = 0.0108$

To compute $\Rightarrow$ probability of $(x/N) \cdot P(N)$

$= 0.0512 \times 0.357$

$= \underline{0.018}$

To compute $p(x/P) \cdot p(P)$

$= 0.0108 \times 0.642$

$= \underline{0.00649}$

$\therefore$ 3 attributes of the rule match with class N.

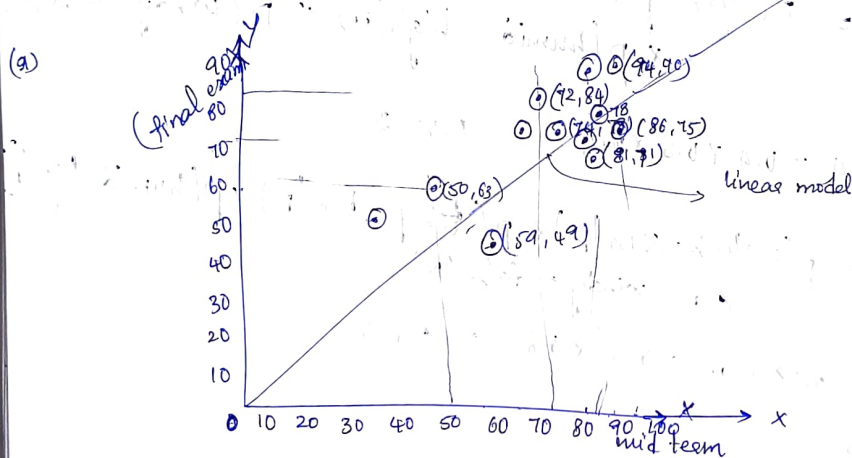=> Univariate Linear model problem (or) Linear regression Problem

① ②

| Mid Team | Final Exam |
|---|---|
| X | Y |
| 72 | 84 |
| 50 | 63 |
| 81 | 71 |
| 74 | 78 |
| 94 | 90 |
| 86 | 90 |
| | 75 |
| 59 | 49 |
| 83 | 79 |
| 65 | 77 |
| 33 | 52 |
| 88 | 74 |
| 81 | 90 |

① Decision tree
   ↓
   ID3 · Algorithm

② Bayesion classification
      ↓
      Naive Bayes

③ linear regression

④ k-means clustering.

a) plot the data x & y
b) use method of test² to find equation for the prediction of student final exam based on the student mid team grades in the grade score:
c) predict final exam grade of student who received 81 marks on the mid team exam.

(a)



X & y are in linear relationship

(b)   $\bar{x}$ = Mean of mid team = $\frac{total}{12}$ = 72.16
      $\bar{y}$ = mean of final exam = $\frac{total}{12}$ = 73.5

(Test square formula)

$$\beta = \frac{\sum_{i=1}^{s} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{s} (x_i - \bar{x})^2}$$

(value for 1st row)

$$\beta = \frac{(72 - 72.16)(84 - 73.5)}{(72 - 72.16)^2} = -65.625$$

$$B_2 = \frac{(50 - 72.16)(63 - 73.5)}{(50 - 72.16)^2} = \frac{(-22.16)(-10.5)}{(22.16)^2} = 0.473$$

$$B_3 = \frac{(81 - 72.16)(71 - 73.5)}{(81 - 72.16)^2}$$

$$B_4 = \frac{(74 - 72.16)(78 - 73.5)}{(74 - 72.16)}$$

$$\boxed{\beta = 0.569} \longrightarrow \text{final answer.}$$

(c)      $\alpha = \bar{y} - \beta \bar{x}$     (predict 87 marks in midterm)

$$= 73.5 - (0.569)(72.16)$$

$$\underline{\underline{\alpha = 32.44}}$$

$$y = \alpha + \beta \cancel{x}$$

$$x \rightarrow 87$$
$$\beta = 0.569$$
$$\alpha = 32.44$$

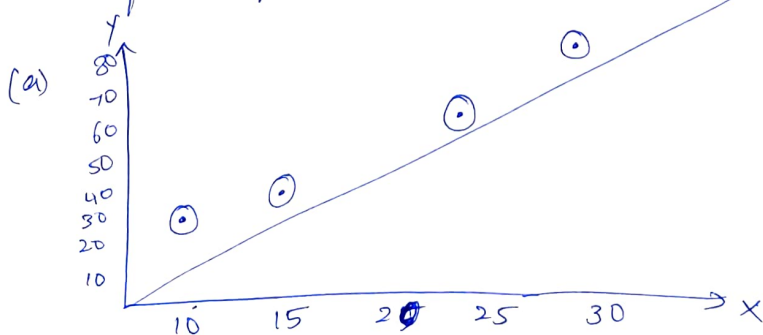$$y = 32.44 + (0.569)(87)$$

$$\boxed{y = 81.943}$$

§ ② 

| X | y |
|---|---|
| 10 | 30 |
| 15 | 40 |
| 25 | 60 |
| 30 | 80 |

(a) same as Q 1
(b) same as Q 2
(c) predict predict 22

(a)

(b)    $\bar{x} = 20$
      $\bar{y} = 52.5$

$\Rightarrow$ k-means clustering problem

Datapoints $(2,3)$, $(6,5)$, $(1,1)$
$(3,3)$, $(8,8)$

% Initial centroids 1 : $(2,3)$
" 2 : $(6,5)$

### step ①:

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

### step ②:

calculate distances from each datapoint to centroid.

1) Distance from $(2,\overset{x_1,y_1}{3})$ to centroid 1 : $(\overset{x_2 \ y_2}{2,3})$

$$d_1 = \sqrt{(x_2-x_1)^2 + (y_2-y_1)^2}$$
$$= \sqrt{(2-2)^2 + (3-3)^2} = \underline{\underline{0}}$$

2) distance from $(3,3)$ to centroid 1 : $(2,3)$

$$d_2 = \sqrt{(3+2)^2 + (3+3)^2} = \sqrt{1^2+0} = \underline{\underline{1}}$$

3) $(6,5)$ to centroid 1 : $(2,3)$

$$d_3 = \sqrt{(6+2)^2 + (5+3)^2} = \sqrt{(4)^2 + 2^2}$$
$$= \sqrt{16+4} = \sqrt{20} = 4.47$$

4) $(8,8)$ to c1 : $(2,3)$

$$d_4 = \sqrt{(6)^2 + (5)^2} = \sqrt{36+25} = \sqrt{61} = 7.8)$$

5) $(1,1)$ to c1 : $(2,3)$

$$d_5 = \sqrt{(1-2)^2 + (1-3)^2} = \sqrt{(1)^2 + (2)^2}$$
$$= \sqrt{1+4} = \sqrt{5} = 2.236$$

Distance to c2 ;

Distance from $(2,3)$ to centroid 2 : $(6,5)$

$$d_1 = \sqrt{4^2 + 2^2} = \sqrt{16+4} = \sqrt{20}$$
$$d_2 = \sqrt{3^2+2^2} = \sqrt{13}$$
$$d_3 = 0$$
$$d_4 = 3.606$$
$$d_5 = 6.403$$

$d1 = 0$ at datapoint $(2,3)$ closest to ∅ centroid 1

$d2 = 1$ at point $(3,3)$ is closest to centroid 1.

$d3 = 0$ at point $(6,5)$ is closest to " 2.

$d4 = 3.606$ at " $(8,8)$ " .. " " 2

$d5 = 2.236$ at " $(1,1)$ " " " (1)

## step③:

Assign points to clusters:

no. of clusters = no- of # centroid.

cluster 1= $(2,3), (3,3)$ $(1,1)$
cluster 2 = $(6,5), (8,8)$

## step④: Recalculate ~~stepwise~~ centroids (iteration 1)

cluster 1: $(2,3)(3,3),(1,1)$

Add $x_1 + x_2, x_3$
Total no- of datapts

new centroid 1 = $\dfrac{x_1 + x_2 + x_3}{3}$ , $\dfrac{y_1 + y_2 + y_3}{3}$

$= \left( \dfrac{2+3+1}{3} \xcancel{} , \dfrac{3+3+1}{3} \right)$

$\begin{array}{r} 6.5 \\ 2\overline{)13} \\ 12 \\ \hline 1 \end{array}$

$= \left( \dfrac{6}{3} , \dfrac{7}{3} \right) = \underline{(2, 2.33)}$

new centroid 2 = $\dfrac{(6,5) , (8,8)}{*_1}$

$= \left( \dfrac{6+8}{2} , \dfrac{5+8}{2} \right)$

$= \left( \dfrac{14}{2} , \dfrac{13}{2} \right) = (7, 6.5)$

$C1, (2,2.3)$ ↓

(d1) $D = \sqrt{(2-2)^2 + (2.3-3)^2} = 0.07$ $(2,3)$ $d1 = \sqrt{(7-2)^2 + (6.5-3)^2} = 6.10$

$C2 (7,6.5)$

(d2) $= \sqrt{(2-3)^2 + (2.3-3)^2} = 1.20$ $(3,3)$ $d2 = \sqrt{(7-3)^2 + (6.5-3)^2} = 5.315$

(d3) $= \sqrt{(2-6)^2 + (2.3-5)^2} = 4.80$ $(6,5)$ (d3) $= \sqrt{(7-6)^2 + (6.5-5)^2} = 1.80$

d4 $= \sqrt{(2-8)^2 + (2.3-8)^2} = 8.24$ $(8,8)$ (d4) $= \sqrt{(7-8)^2 + (6.5-8)^2} = 1.80$

(d5) $= \sqrt{(2-1)^2 + (2.3-1)^2} = 1.66$ $(1,1)$ $d5 = \sqrt{(7-1)^2 + (6.5-1)^2} = 8.139$

new cluster assignment:
cluster 1 = $(2,3) (3,3) (1,1)$
cluster 2 = $( 6,5), (8,8)$

clusters do not change after the 2nd iteration

**final clusters**

C1: (2,3) (3,3), (1,1)

C2: (6,5) (6,8)    ; int iteration 2

§ (4,7), (5,6), (7,7), (2,3),. (6,5)

Centroid c1: (4,7)

centroid c2: (7,7)

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

Centroid c1;

(d1) wit = $\sqrt{0+0}$  = 0

(d2) = $\sqrt{(-1)^2+(1)^2}$ = $\sqrt{1+1}$ = $\sqrt{2}$ = 1.414

d3 = $\sqrt{(-3)^2+(0)^2}$ = $\sqrt{9}$ = 3 *

(d4) = $\sqrt{2^2+4^2}$ = $\sqrt{4+16}$ = $\sqrt{20}$ = 4.472

d5 = $\sqrt{1^2+2^2}$ = $\sqrt{1+4}$ = $\sqrt{5}$ = 2.236

$\sqrt{(-2)^2+2^2}$ = $\sqrt{4+4}$ : $\sqrt{8}$ = 2.828

centroid c2;

d1 = $\sqrt{3^2+0^2}$ = $\sqrt{9}$ = 3

d2 = $\sqrt{2^2+1^2}$ = $\sqrt{4+1}$ = $\sqrt{5}$ = 2.236

(d3) = $\sqrt{0+0}$ = 0

d4 = $\sqrt{5^2+4^2}$ = $\sqrt{25+16}$ = $\sqrt{41}$ = 6.403

(d5) = $\sqrt{1^2+2^2}$ = $\sqrt{5}$ = 2.236

cluster 1: (4,7), (5,6), (2,3)

cluster 2: (7,7), (6,5)

new centroid for cluster 1: $\left(\frac{4+5+2}{3}, \frac{7+6+3}{3}\right)$

= $\left(\frac{11}{3}, \frac{16}{3}\right)$ = (3.66, 5.33)

new centroid for cluster 2: $\left(\frac{7+6}{2}, \frac{7+5}{2}\right)$

= $\left(\frac{13}{2}, \frac{12}{2}\right)$

= (6.5, 6)

new centroid c1;

$d1 = \sqrt{0.1156 + 2.7889} = 1.704$

$d2 = \sqrt{1.79 + 0.4489} = 1.49$

$d3 = \sqrt{11.15 + 2.78} = 3.73$

$d4 = \sqrt{2.15 + 5.42} = 2.85$

$d5 = \sqrt{5.47 + 0.10} = 2.36$

new centroid c2;

$d1 = \sqrt{6.25 + 1} = 2.69$

$d2 = \sqrt{2.25 + 0} = 5.06$

$d3 = \sqrt{0.25 + 1} = 1.11$

$d4 = \sqrt{20.25 + 9} = \sqrt{29.25} = 5.40$

$d5 = \sqrt{0.25 + 1} = 1.11$

cluster 1 : (4,7) (5,6) (7,7) 2,3)
cluster 2: {(7,7) (6,5)

new centroid c1 => (3.66, 5.33)
new centroid c2 => (6.5, 6)

∴ Iteration 2 will also give the same.