

CS 6350 - BIG DATA ANALYTICS & MGT  
HOMEWORK - 4

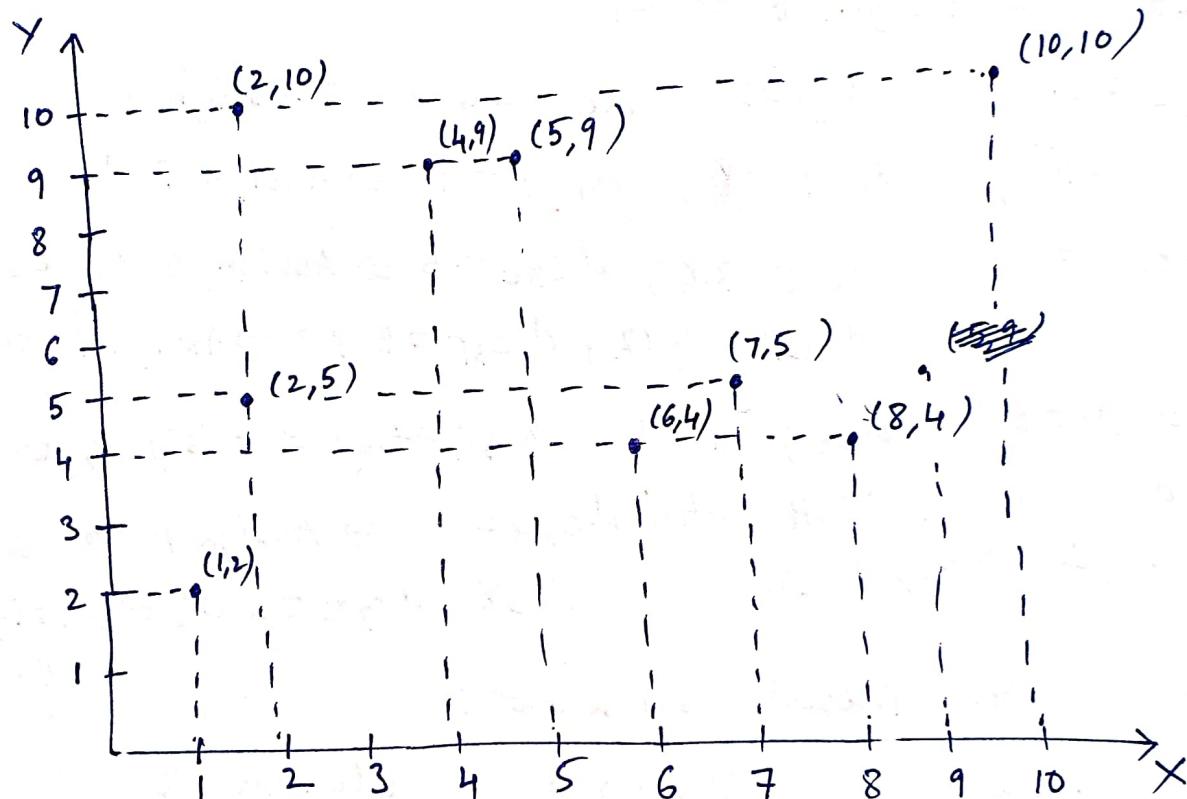
PART-1 : CLUSTERING .

QUESTION A :

Points :  $(2,10), (2,5), (8,4), (5,9), (7,5), (6,4), (1,2), (4,9), (10,10)$  .

Let's label these points from A to I

- (i). Plot the data



Possible clusters :

Cluster-1

$(2,5)$

$(1,2)$

Cluster-2 :

$(2,10)$

$(4,9)$

$(5,9)$

Cluster-3

$(7,5)$

~~$(8,4)$~~

$(6,4)$

$(8,4)$

$(10,10)$

(ii). Initial cluster centres .

$$C_1 = (2, 5)$$

$$C_2 = (5, 8)$$

$$C_3 = (4, 9)$$

Assign points to nearest clusters.

$$d_{C_1 A} = 5, d_{C_2 A} = 3.60, d_{C_3 A} = 2.24 \Rightarrow \text{Assign } A \text{ to } C_3$$

$$d_{C_1 B} = 0 \Rightarrow \text{Assign } B \text{ to } C_1$$

$$d_{C_1 C} = 6.08, d_{C_2 C} = 5, d_{C_3 C} = 6.4 \Rightarrow \text{Assign } C \text{ to } C_2$$

$$d_{C_1 D} = 5, d_{C_2 D} = 10, d_{C_3 D} = 1 \Rightarrow \text{Assign } D \text{ to } C_2$$

$$d_{C_1 E} = 5, d_{C_2 E} = 3.6, d_{C_3 E} = 5 \Rightarrow \text{Assign } E \text{ to } C_2$$

$$d_{C_1 F} = 4.12, d_{C_2 F} = 4.12, d_{C_3 F} = 5.38 \Rightarrow \text{Assign } F \text{ to } C_2$$

$$d_{C_1 G} = 3.16, d_{C_2 G} = 7.21, d_{C_3 G} = 7.61 \Rightarrow \text{Assign } G \text{ to } C_1$$

$$d_{C_1 H} = 4.47, d_{C_2 H} = 1.41, d_{C_3 H} = 0 \Rightarrow \text{Assign } H \text{ to } C_3$$

$$d_{C_1 I} = 9.43, d_{C_2 I} = 5.38, d_{C_3 I} = 6.08 \Rightarrow \text{Assign } I \text{ to } C_2$$

Q Clusters after initial assignment :

Cluster 1

$$B = (2, 5)$$

~~$$C = (8, 4)$$~~

$$G = (1, 2)$$

Cluster 2

$$D = (5, 8)$$

$$E = (7, 5)$$

$$F = (6, 4)$$

$$I = (10, 10)$$

$$C = (8, 4)$$

Cluster 3

~~$$A = (2, 10)$$~~

$$H = (4, 9)$$

## Round-1 (contd)

Calculating the centroids of the clusters

$$c_1 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

$$c_2 = \left( \frac{5+7+6+10+8}{5}, \frac{9+5+4+10+4}{5} \right) = (7.2, 6.4)$$

$$c_3 = \left( \frac{2+4}{2}, \frac{10+9}{2} \right) = (3, 9.5)$$

## Round-2 :

Assign points to closest centroids

$$dc_{1A} = 6.52, dc_{2A} = 6.32, dc_{3A} = 1.12 \Rightarrow \text{Assign A to } c_3$$

$$dc_{1B} = 1.58, dc_{2B} = 5.38, dc_{3B} = 4.6 \Rightarrow \text{Assign B to } c_1$$

$$dc_{1C} = 6.52, dc_{2C} = 2.53, dc_{3C} = 7.43 \Rightarrow \text{Assign C to } c_2$$

$$dc_{1D} = 6.52, dc_{2D} = 3.4, dc_{3D} = 2.06 \Rightarrow \text{Assign D to } c_3$$

$$dc_{1E} = 5.7, dc_{2E} = 1.41, dc_{3E} = 6.02 \Rightarrow \text{Assign E to } c_2$$

$$dc_{1F} = 4.53, dc_{2F} = 2.68, dc_{3F} = 6.26 \Rightarrow \text{Assign F to } c_2$$

$$dc_{1G} = 1.58, dc_{2G} = 7.60, dc_{3G} = 7.76 \Rightarrow \text{Assign G to } c_1$$

$$dc_{1H} = 6.04, dc_{2H} = 4.20, dc_{3H} = 1.12 \Rightarrow \text{Assign H to } c_3$$

$$dc_{1I} = 10.7, dc_{2I} = 4.56, dc_{3I} = 7.02 \Rightarrow \text{Assign I to } c_2$$

Cluster 1      Cluster 2      Cluster 3

$$B = (2, 5)$$

$$C = (8, 4)$$

$$H = (4, 9)$$

$$G = (1, 2)$$

$$E = (7, 5)$$

$$D = (5, 9)$$

$$F = (6, 4)$$

$$A = (2, 10)$$

$$I = (10, 10)$$

Round-2 (contd.)

$$C_1 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

$$C_2 = \left( \frac{8+7+6+9+10}{4}, \frac{4+5+4+10}{4} \right) = (7.75, 5.75)$$

$$C_3 = \left( \frac{4+5+2}{3}, \frac{9+9+10}{3} \right) = (3.67, 9.33)$$

Round 3:

$$d_{C_1 A} = 6.52, d_{C_2 A} = 7.15, d_{C_3 A} = 1.79 \Rightarrow \text{Assign } A \text{ to } C_3$$

$$d_{C_1 B} = 1.58, d_{C_2 B} = 5.79, d_{C_3 B} = 4.64 \Rightarrow \text{Assign } B \text{ to } C_1$$

$$d_{C_1 C} = 6.52, d_{C_2 C} = 1.76, d_{C_3 C} = 6.86 \Rightarrow \text{Assign } C \text{ to } C_2$$

$$d_{C_1 D} = 6.52, d_{C_2 D} = \cancel{2.23} 4.25, d_{C_3 D} = 1.37 \Rightarrow \text{Assign } D \text{ to } C_3$$

$$d_{C_1 E} = 5.7, d_{C_2 E} = 1.06, d_{C_3 E} = 5.46 \Rightarrow \text{Assign } E \text{ to } C_2$$

$$d_{C_1 F} = 4.53, d_{C_2 F} = 2.47, d_{C_3 F} = 5.81 \Rightarrow \text{Assign } F \text{ to } C_2$$

$$d_{C_1 G} = 1.58, d_{C_2 G} = 7.72, d_{C_3 G} = 7.8 \Rightarrow \text{Assign } G \text{ to } C_1$$

$$d_{C_1 H} = 6.03, d_{C_2 H} = 4.96, d_{C_3 H} = 0.46 \Rightarrow \text{Assign } H \text{ to } C_3$$

$$d_{C_1 I} = 10.7, d_{C_2 I} = 4.8, d_{C_3 I} = 6.36 \Rightarrow \text{Assign } I \text{ to } C_2$$

Cluster 1

$$B = (2, 5)$$

$$G = (1, 2)$$

Cluster 2

$$C = (8, 4)$$

$$E = (7, 5)$$

$$F = (6, 4)$$

$$I = (10, 10)$$

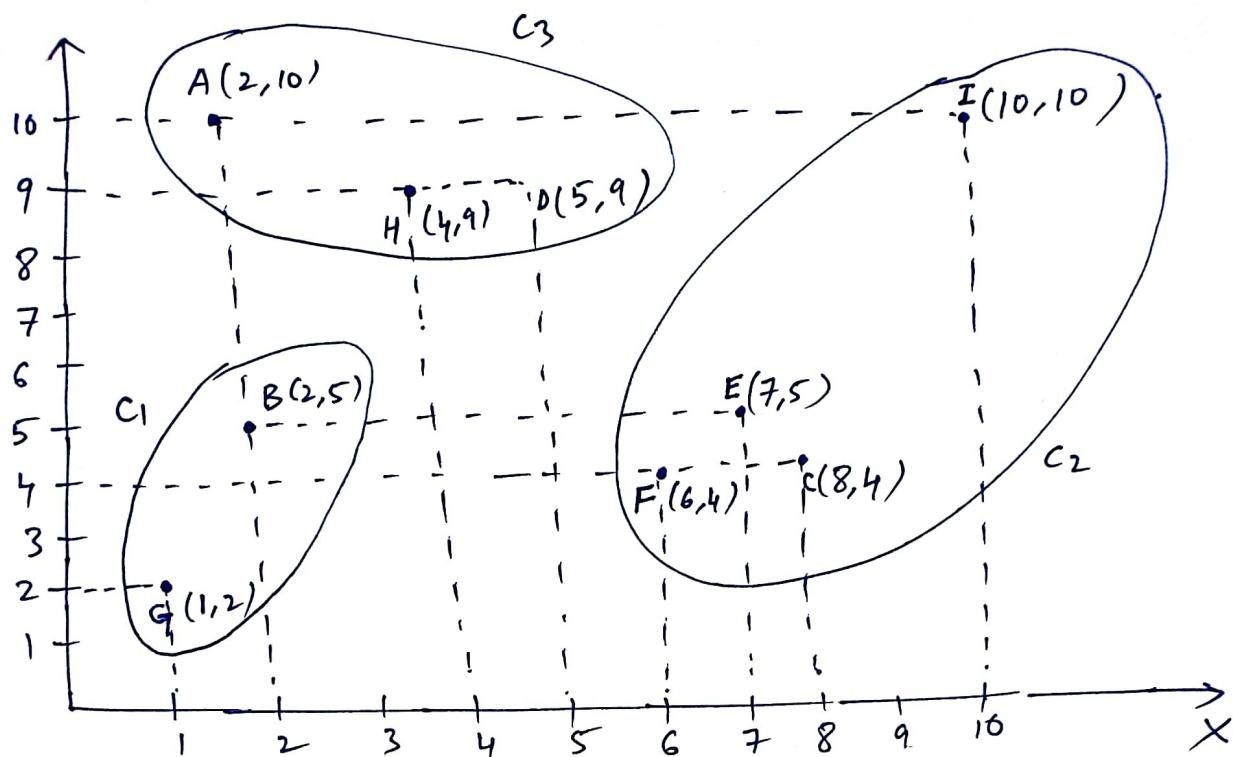
Cluster 3

$$H = (4, 9)$$

$$D = (5, 9)$$

$$A = (2, 10)$$

Since the clusters don't change in round 3,  
this is the final set of clusters.



Cluster centres :

For C<sub>1</sub> : (1.5, 3.5)

For C<sub>2</sub> : (7.75, 5.75)

for C<sub>3</sub> : (3.67, 9.33)

### QUESTION B :

(i). Single link hierarchical clustering

#### Round-1:

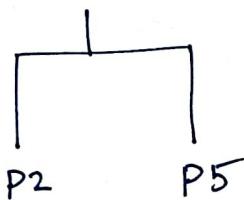
Clusters  $P_2$  and  $P_5$  are most similar with a value of 0.98.

Initial

Clusters :

$\{P_1\}, \{P_2\}, \{P_3\},$   
 $\{P_4\}, \{P_5\}$

Dendrogram :



Updating the proximity matrix .

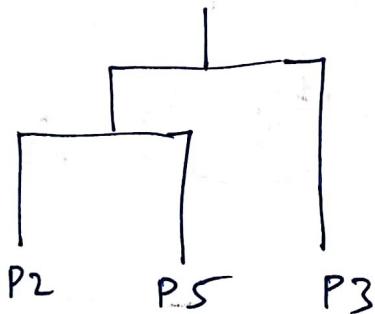
	$(P_2, P_5)$	$P_1$	$P_3$	$P_4$
$(P_2, P_5)$	1	0.35	0.85	0.76
$P_1$	0.35	1	0.41	0.55
$P_3$	0.85	0.41	1	0.44
$P_4$	0.76	0.55	0.44	1

$$\text{eg. } \text{sim}\{(P_2, P_5), P_1\} \\ = \max\{\text{sim}(P_2, P_1), \\ \text{sim}(P_5, P_1)\} \\ = \max(0.1, 0.35) \\ = 0.35.$$

#### Round-2:

Clusters  $\{P_2, P_5\}$  and  $\{P_3\}$  are most similar with a value of 0.85, and so we merge them .

Dendrogram:

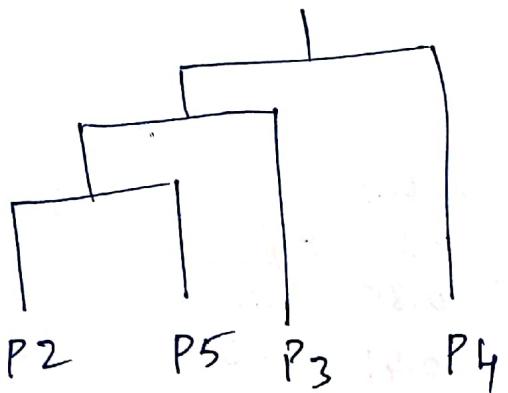


### Updated proximity matrix

	(P <sub>2</sub> , P <sub>5</sub> , P <sub>3</sub> )	P <sub>1</sub>	P <sub>4</sub>
(P <sub>2</sub> , P <sub>5</sub> , P <sub>3</sub> )	1	0.41	0.76
P <sub>1</sub>	0.41	1	0.55
P <sub>4</sub>	0.76	0.55	1

Round-3:

Clusters {P<sub>2</sub>, P<sub>5</sub>, P<sub>3</sub>} and {P<sub>4</sub>} are most similar with a value of 0.76, hence we merge them.

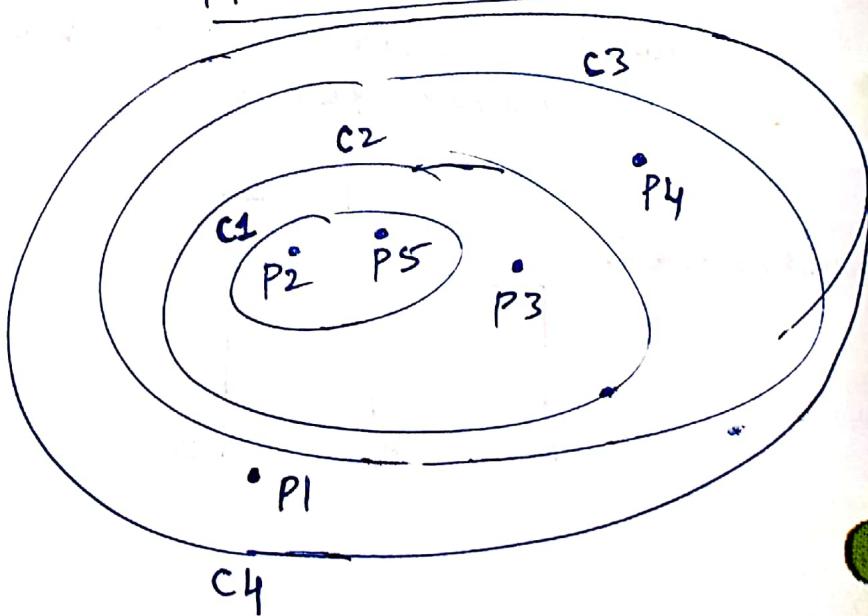
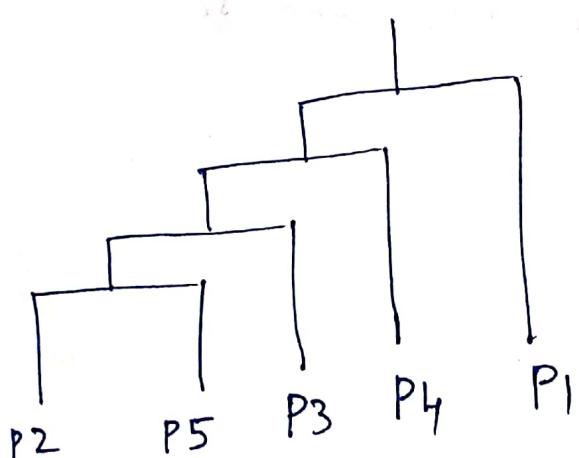


~~Round~~

Updated similarity matrix is not required -

Since only 2 clusters are left, we merge them.

Final hierarchical clustering

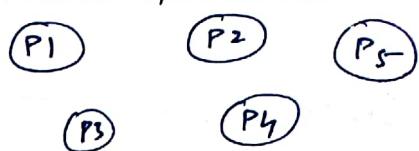


## Complete link hierarchical clustering:

### Round-1:

Clusters  $\{P_2\}$  and  $\{P_5\}$  have the most similarity ( $\text{value} = 0.98$ ), so we merge them.

### Initial clusters:



### Dendrogram:



### Updated similarity matrix:

	(P2, P5)	P1	P3	P4
(P2, P5)	1	0.1	0.64	0.47
P1	0.1	1	0.41	0.55
P3	0.64	0.41	1	0.44
P4	0.47	0.55	0.44	1

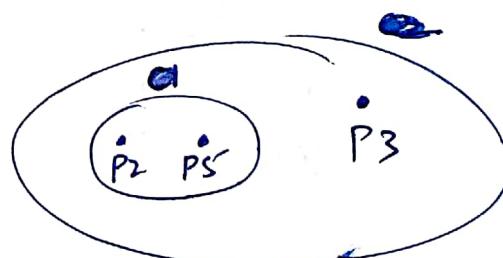
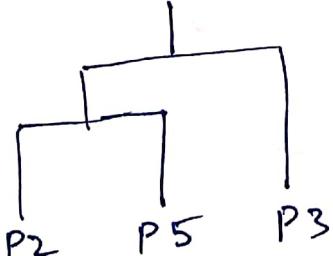
### Example:

$$\begin{aligned} & \text{sim}\{(P_2, P_5), P_1\} \\ &= \min\{\text{sim}(P_2, P_1), \\ & \quad \text{sim}(P_5, P_1)\} \\ &= \min(0.1, 0.35) \\ &= 0.1 \end{aligned}$$

### Round-2:

$(P_2, P_5)$  and  $P_3$  have the maximum similarity value = 0.64, so we merge them.

### Dendrogram:



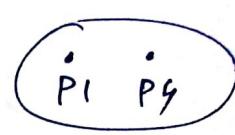
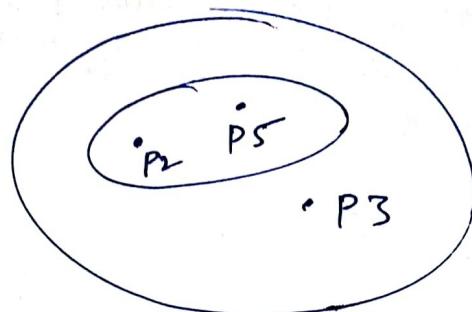
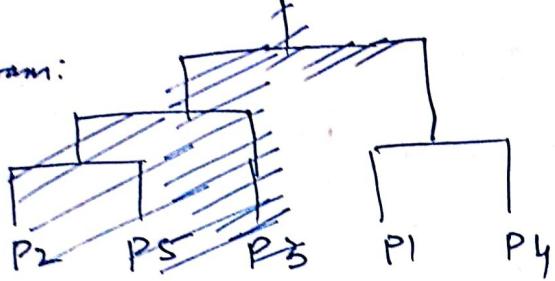
### Updated similarity matrix:

	(P2, P5, P3)	P1	P4
(P2, P5, P3)	1	0.1	0.44
P1	0.1	1	0.55
P4	0.44	0.55	1

### Round-3:

Clusters  $\{P_1\}$  and  $\{P_4\}$  have the highest similarity value of 0.55, so merge them.

Dendrogram:



Updated Similarity matrix: Not required (since only 2 clusters left)

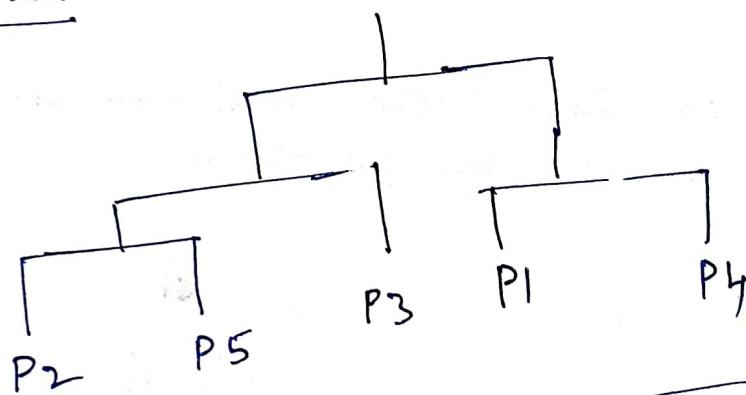
### Round-4:

Since only 2 clusters remain:

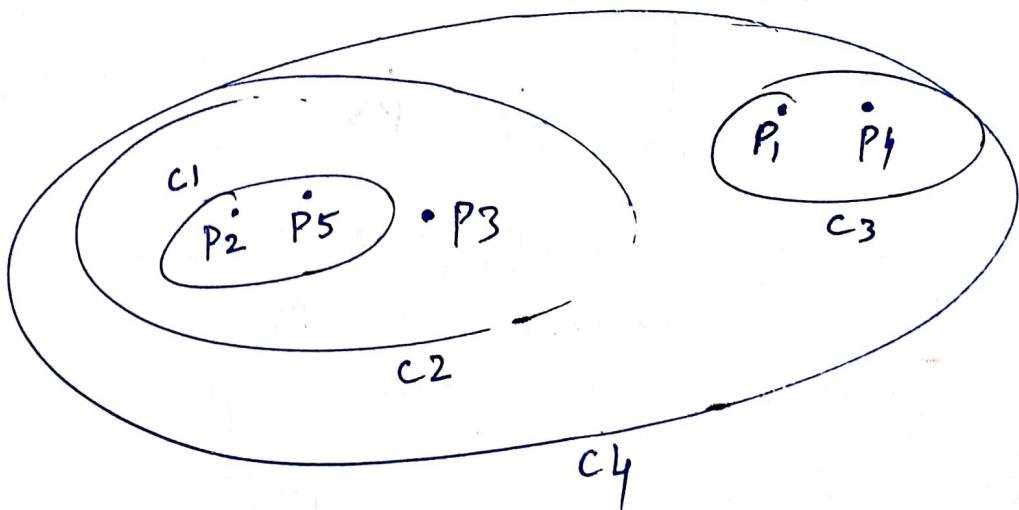
$\{(P_2, P_5), P_3\}$  and  $\{P_1, P_4\}$

we merge them.

Final Dendrogram:



Final hierarchical cluster:



QUESTION C:

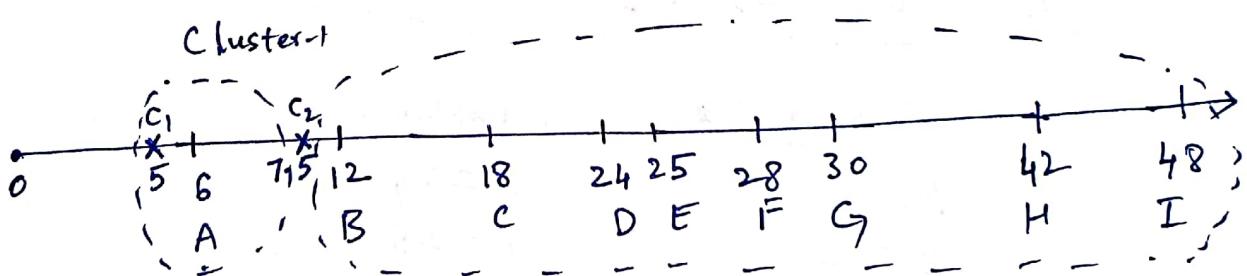
(a). Set of points:  $A=6, B=12, C=18, D=24$   
 $E=25, F=28, G=30, H=42,$   
 $I=48$ .

(a). Set of centroids:

$$c_1 = 5$$

$$c_2 = 7.5$$

$d_{C_1 A} = 1, d_{C_2 A} = 1.5 \Rightarrow$  Assign A to  $c_1$   
 $d_{C_1 B} = 7, d_{C_2 B} = 4.5 \Rightarrow$  Assign B to  $c_2$   
 $d_{C_1 C} = 13, d_{C_2 C} = 10.5 \Rightarrow$  Assign C to  $c_2$   
 $d_{C_1 D} = 19, d_{C_2 D} = 16.5 \Rightarrow$  Assign D to  $c_2$   
 $d_{C_1 E} = 20, d_{C_2 E} = 17.5 \Rightarrow$  Assign E to  $c_2$   
 $d_{C_1 F} = 23, d_{C_2 F} = 20.5 \Rightarrow$  Assign F to  $c_2$   
 $d_{C_1 G} = 25, d_{C_2 G} = 22.5 \Rightarrow$  Assign G to  $c_2$   
 $d_{C_1 H} = 37, d_{C_2 H} = 34.5 \Rightarrow$  Assign H to  $c_2$   
 $d_{C_1 I} = 43, d_{C_2 I} = 40.5 \Rightarrow$  Assign I to  $c_2$



Cluster 1      Cluster 2

$$A = 6$$

$$B = 12$$

$$C = 18$$

$$D = 24$$

$$E = 25$$

$$F = 28$$

$$G = 30$$

$$H = 42$$

$$I = 48$$

Total squared error.

$$\begin{aligned}
 &= (1)^2 + (4.5)^2 + (10.5)^2 \\
 &\quad + (16.5)^2 + (17.5)^2 \\
 &\quad + (20.5)^2 + (22.5)^2 \\
 &\quad + (34.5)^2 + (40.5)^2
 \end{aligned}$$

$$= 4467$$

(2). Set of centroids:

$$C_1 = 15$$

$$C_2 = 25$$

$d_{C_1 A} = 9, d_{C_2 A} = 19 \Rightarrow$  Assign A to  $C_1$

$d_{C_1 B} = 3, d_{C_2 B} = 13 \Rightarrow$  assign B to  $C_1$

$d_{C_1 C} = 3, d_{C_2 C} = 7 \Rightarrow$  assign C to  $C_1$

$d_{C_1 D} = 9, d_{C_2 D} = 1 \Rightarrow$  assign D to  $C_2$

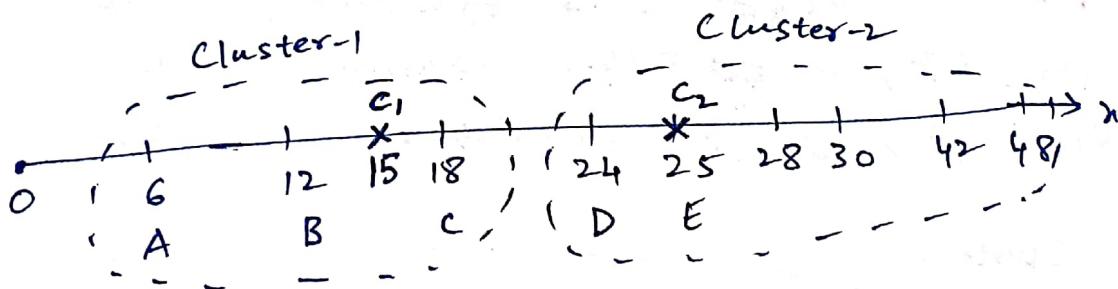
$d_{C_1 E} = 10, d_{C_2 E} = 0 \Rightarrow$  assign E to  $C_2$

$d_{C_1 F} = 13, d_{C_2 F} = 3 \Rightarrow$  assign F to  $C_2$

$d_{C_1 G} = 15, d_{C_2 G} = 5 \Rightarrow$  assign G to  $C_2$

$d_{C_1 H} = 27, d_{C_2 H} = 17 \Rightarrow$  assign H to  $C_2$

$d_{C_1 I} = 33, d_{C_2 I} = 23 \Rightarrow$  assign I to  $C_2$



Cluster-1

$$A = 6$$

$$B = 12$$

$$C = 18$$

Cluster-2

$$D = 24$$

$$E = 25$$

$$F = 28$$

$$G = 30$$

$$H = 42$$

$$I = 48$$

Total sum of errors =  $9^2 + 3^2 + 3^2 + 1^2 + 0^2 + 3^2 + 5^2 + 17^2 + 23^2$   
 $= 898$

(b). Centroid set 1

$$c_1 = 5$$

$$c_2 = 7.5.$$

Assume the points were allocated as in part a.  
Recalculating the new centroids

$$c_1 = 6$$

$$c_2 = \frac{12+18+24+25+28+30+42+48}{8} = 28.37.$$

$d_{c_1 A} = 0$ ,  $d_{c_2 A} = 22.37 \Rightarrow$  assign A to  $c_1$

$d_{c_1 B} = 6$ ,  $d_{c_2 B} = 16.37 \Rightarrow$  assign B to  $c_1$

$d_{c_1 C} = 12$ ,  $d_{c_2 C} = 10.37 \Rightarrow$  assign C to  $c_2$

$d_{c_1 D} = 18$ ,  $d_{c_2 D} = 4.37 \Rightarrow$  assign D to  $c_2$

$d_{c_1 E} = 19$ ,  $d_{c_2 E} = 13.37 \Rightarrow$  assign E to  $c_2$

$d_{c_1 F} = 22$ ,  $d_{c_2 F} = 0.37 \Rightarrow$  assign F to  $c_2$

$d_{c_1 G} = 24$ ,  $d_{c_2 G} = 1.63 \Rightarrow$  assign G to  $c_2$

$d_{c_1 H} = 36$ ,  $d_{c_2 H} = 13.63 \Rightarrow$  assign H to  $c_2$

$d_{c_1 I} = 42$ ,  $d_{c_2 I} = 19.63 \Rightarrow$  assign I to  $c_2$

Cluster-1:

$$A = 6$$

$$B = 12$$

Cluster 2

$$C = 18$$

$$D = 24$$

$$E = 25$$

$$F = 28$$

$$G = 30$$

$$H = 42$$

$$I = 48$$

Since there is a change in points in each cluster, centroid set  $= (5, 7.5)$  is not a stable solution.

### Centroid set - 2 :

Assume that the points were allocated to the 2 clusters as in part (a)  
Calculating the centroids (updated)

$$C_1 = \left( \frac{6+12+18}{3} \right), \left( \frac{24+25+28+30+42+48}{6} \right) \\ = (12, 32.83)$$

$$C_1 = 12$$

$$C_2 = 32.83$$

Round 2: Assign points to the 2 clusters.

$d_{C_1 A} = 6, d_{C_2 A} = 26.83 \Rightarrow$  assign A to  $C_1$

$d_{C_1 B} = 0, d_{C_2 B} = 20.83 \Rightarrow$  assign B to  $C_1$

$d_{C_1 C} = 6, d_{C_2 C} = 14.83 \Rightarrow$  assign C to  $C_1$

$d_{C_1 D} = 12, d_{C_2 D} = 8.83 \Rightarrow$  assign D to  $C_2$

$d_{C_1 E} = 13, d_{C_2 E} = 7.83 \Rightarrow$  assign E to  $C_2$

$d_{C_1 F} = 16, d_{C_2 F} = 4.83 \Rightarrow$  assign F to  $C_2$

$d_{C_1 G} = 18, d_{C_2 G} = 2.83 \Rightarrow$  assign G to  $C_2$

$d_{C_1 H} = 30, d_{C_2 H} = 9.17 \Rightarrow$  assign H to  $C_2$

$d_{C_1 I} = 36, d_{C_2 I} = 15.17 \Rightarrow$  assign I to  $C_2$

### Cluster-1      Cluster-2

$$A = 6 \qquad D = 24$$

$$B = 12 \qquad E = 25$$

$$C = 18 \qquad F = 28$$

$$G = 30$$

$$H = 42$$

$$I = 48$$

∴ Centroid set = (15, 25) represents a stable solution.

(c). MIN/single link clustering.

Proximity matrix:

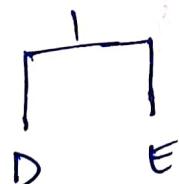
	A	B	C	D	E	F	G	H	I
A	0	6	12	18	19	22	24	36	42
B	6	0	6	12	13	16	18	30	36
C	12	6	0	6	7	10	12	24	30
D	18	12	6	0	1	4	6	18	24
E	19	13	7	1	0	3	5	17	23
F	22	16	10	4	3	0	2	14	20
G	24	18	12	6	5	2	0	12	16
H	36	30	24	18	17	14	12	0	6
I	42	36	30	24	23	20	16	6	0

Initially:



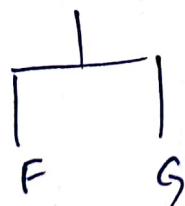
Round 1:

Clusters containing D and E are closest, so merge them ( $DE=1$ )



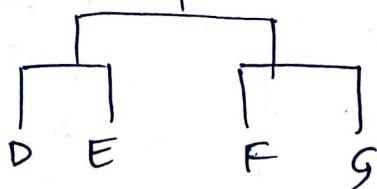
Round-2:

Clusters {F} and {G} are closest so merge them.  
( $FG=2$ )



### Round-3 :

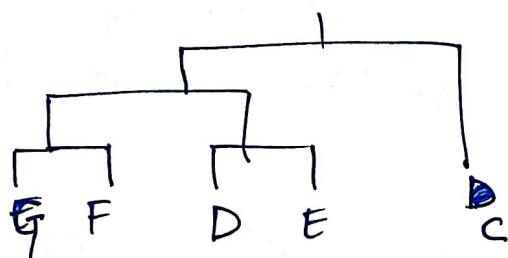
Clusters  $\{D, E\}$  and  $\{F, G\}$  are closest, so merge them ( $FE=3$ )



(~~D/E/F/G~~)

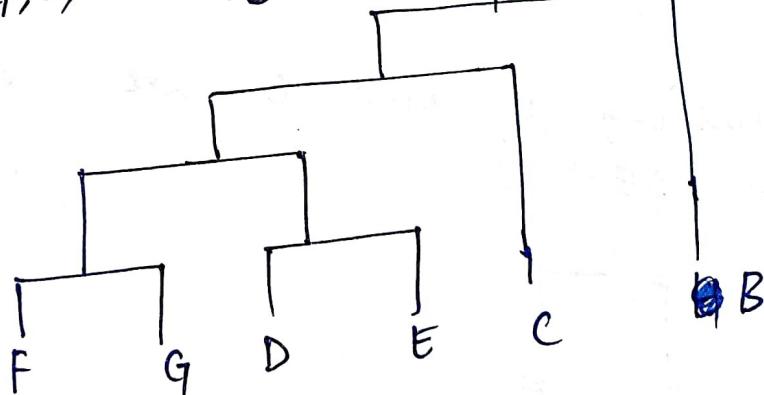
### Round-4 :

Clusters  $\{C\}$  and  $\{D, E, F, G\}$  are closest, so we merge them ( $CE=7$ )



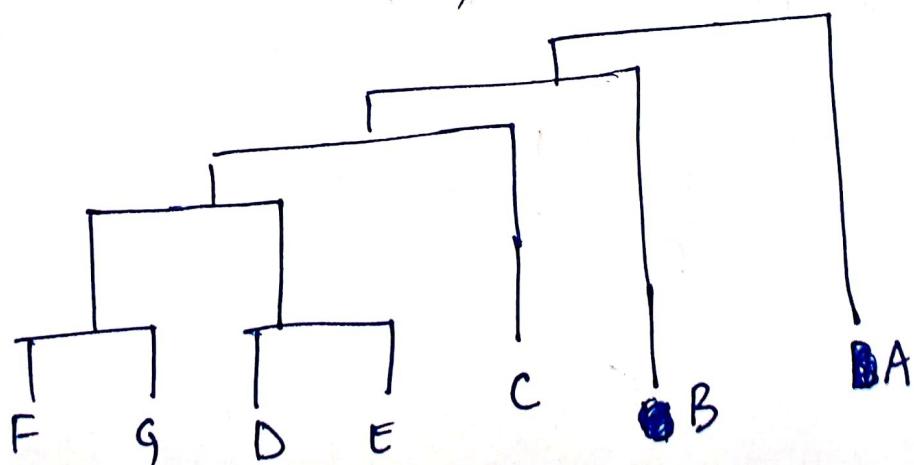
### Round-5 :

Merge  $\{E, F, D, G, C\}$  and  $\{B\}$  (~~BC=6~~)



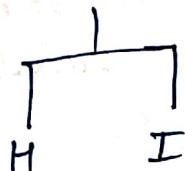
### Round-6 :

Merge  $\{F, G, D, E, C, H\}$  and  $\{A\}$  (~~AB=6~~)



### Round-7:

We Merge  $\{H\}$  and  $\{I\}$  since  $HI=6$  is the lowest distance between any 2 clusters.



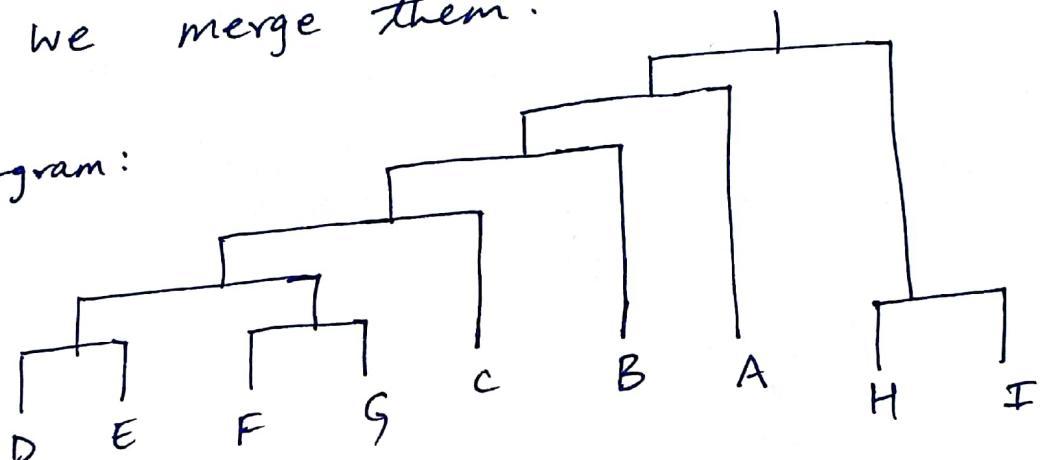
### Round-8:

Since only 2 clusters remain:

$\{A, B, C, D, E, F, G\}$  and  $\{H, I\}$

we merge them.

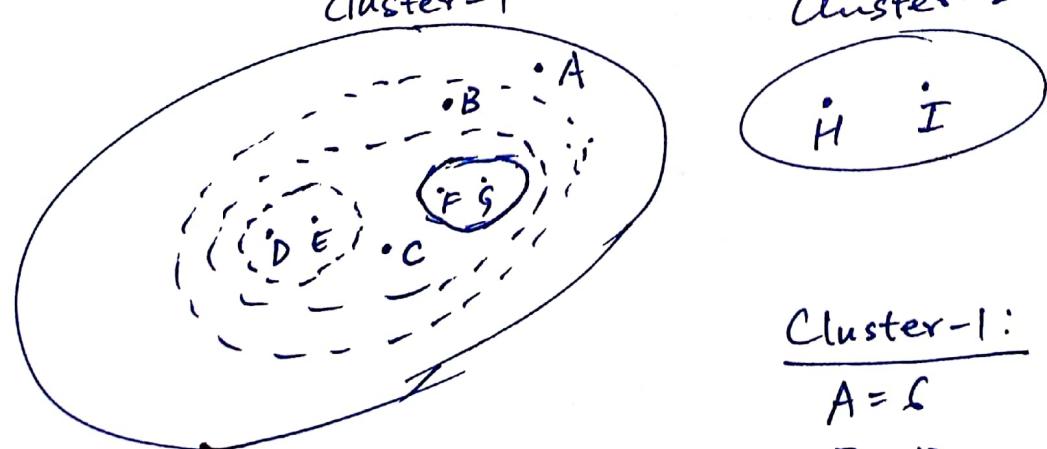
Final dendrogram:



If we want 2 clusters, we must go down one level in the final dendrogram

Cluster-1

Cluster-2



Cluster-1:

$$A = 5$$

$$B = 12$$

$$C = 18$$

$$D = 24$$

$$E = 25$$

$$F = 28$$

$$G = 30$$

Cluster-2:

$$H = 42$$

$$I = 48$$

### K-means

- (d). ~~PEEL~~ seems to produce the most natural clustering.
- (e). The K-means algorithm has the advantage of centroid updation for each cluster in each round, which helps in similar points getting merged into the same cluster.

## PART II : CLASSIFICATION

### QUESTION D :

Label the examples from 1 to 11.

$$S_0 = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11\}$$

Target attribute :  $y$

$$x_1 \in \{a, b\}, x_2 \in \{c, g, u, w\}, x_3 \in \{k, s, v\}$$

$x$  creates a partition of (5,6)

$$\text{Entropy}(S_0) = \text{Entropy}\left(\frac{5}{11}, \frac{6}{11}\right) = 0.517 + 0.476 = 0.9939$$

### Selecting attribute

$x_1 = a$  creates a partition of (2,3)

$$\text{Entropy}(S_0|x_1=a) = \text{Entropy}\left(\frac{2}{5}, \frac{3}{5}\right) = 0.823.$$

$x_1 = b$  creates a partition of (3,3)

$$\text{Entropy}(S_0|x_1=b) = \text{Entropy}\left(\frac{3}{6}, \frac{3}{6}\right) = 1$$

$$\therefore \text{Entropy}(S_0|x_1) = \frac{5}{11}(0.823) + \frac{6}{11}(1) = 0.9195 \quad \text{--- (1)}$$

$x_2 = c$  creates a partition of (3,2)

$$\text{Entropy}(S_0|x_2=c) = 0.823$$

$x_2 = g$  creates a partition of (1,0)

$$\text{Entropy}(S_0|x_2=g) = \text{Entropy}(1,0) = 0$$

$x_2 = u$  creates a partition of (0,1)

$$\text{Entropy}(S_0|x_2=u) = \text{Entropy}(0,1) = 0$$

$x_2 = w$  creates a partition of (2,2)

$$\text{Entropy}(S_0|x_2=w) = 1$$

$$\begin{aligned} \therefore \text{Entropy}(S_0|x_2) &= \frac{5}{11} \times (0.823) + \frac{1}{11}(0) + \frac{1}{11}(0) + \frac{4}{11}(1) \\ &= 0.7377 \quad \text{--- (2)} \end{aligned}$$

$x_3 = K$  creates a partition of  $(0, 3)$

$$\text{Entropy}(S_0 | x_3 = K) = \text{Ent}(0, 1) = 0$$

$x_3 = V$  creates a partition of  $(2, 3)$

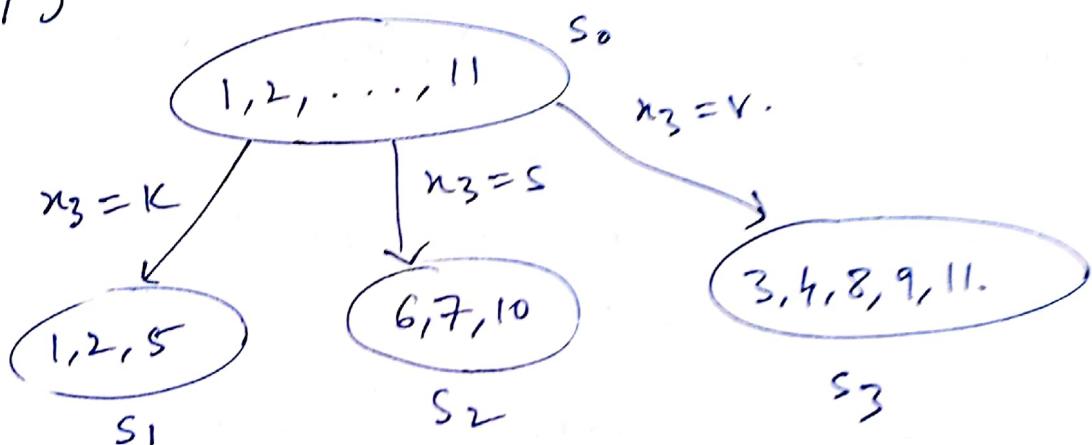
$$\text{Entropy}(S_0 | x_3 = V) = \text{Entropy}\left(\frac{2}{5}, \frac{3}{5}\right) = 0.823$$

$x_3 = S$  creates a partition  $(3, 0)$

$$\text{Entropy}(S_0 | x_3 = S) = \text{Entropy}(1, 0) = 0$$

$$\therefore \text{Entropy}(S_0 | x_3) = \frac{3}{11}(0) + \frac{5}{11}(0.823) + \frac{3}{11}(0) \\ = 0.374.$$

We split based on the attribute having lowest entropy.



All examples in  $S_1$  have same  $Y$ .

All examples in  $S_2$  have same  $Y$

$$S_3 = \{3, 4, 8, 9, 11\}$$

$x_1 = a$  creates a partition of  $(1, 1)$

$$\text{Entropy}(S_3 | x_1 = a) = 1$$

$x_1 = b$  creates a partition of  $(1, 2)$

$$\text{Entropy}(S_3 | x_1 = b) = 0.9179$$

$$\text{Entropy}(S_3 | x_1) = \frac{2}{5}(1) + \frac{3}{5}(0.9179) = 0.9507 \quad \text{--- (3)}$$

$x_2 = c$  creates a  $(1, 1)$  partition

$$\text{Entropy}(S_3 | x_2 = c) = \text{Entropy}(\frac{1}{2}, \frac{1}{2}) = 1$$

$x_2 = g$  creates a partition  $(0, 1)$

$$\text{Entropy}(S_3 | x_2 = g) = 0$$

$x_2 = u$  creates partition  $(0, 0)$

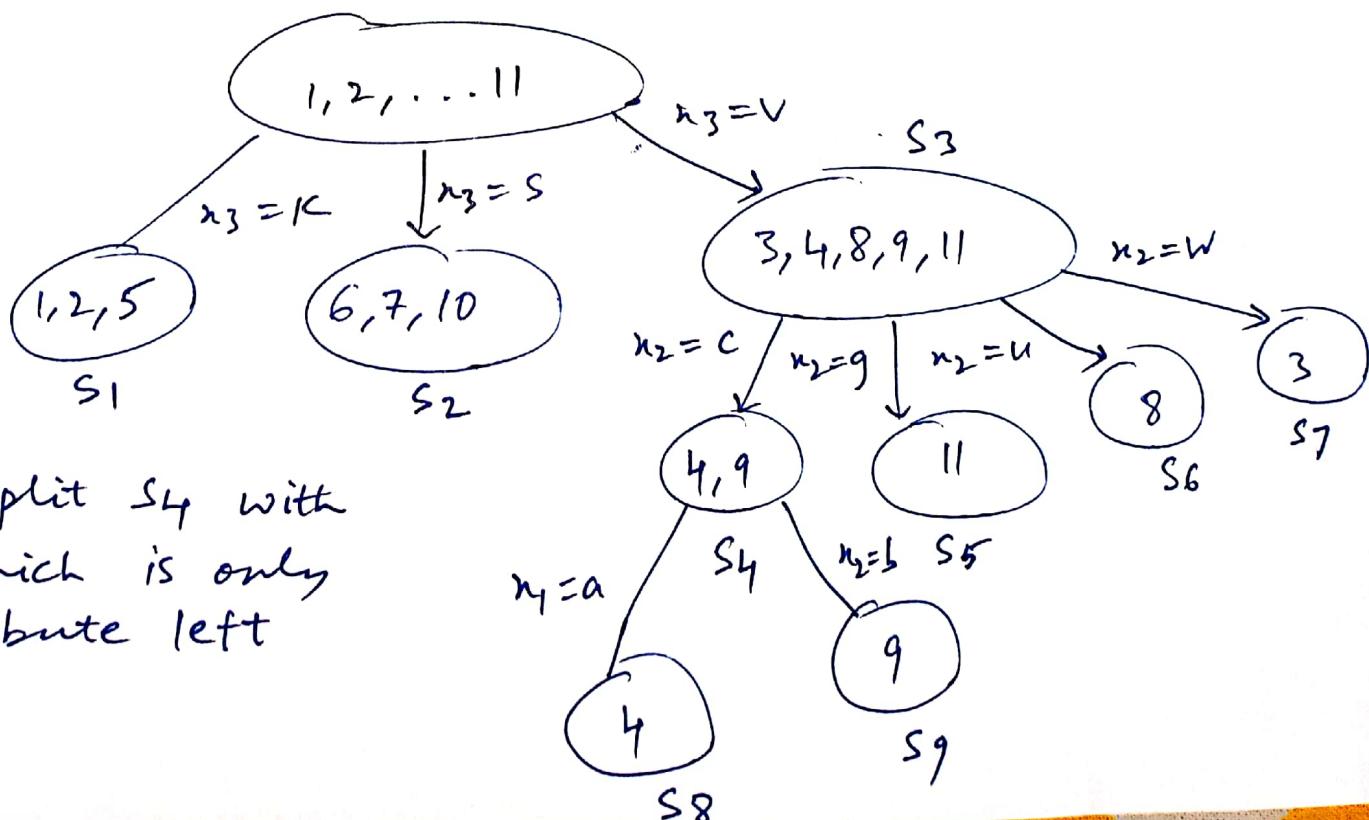
$$\text{Entropy}(S_3 | x_2 = u) = 0$$

$x_2 = w$  creates a  $(1, 0)$  partition

$$\text{Entropy}(S_3 | x_2 = w) = 0$$

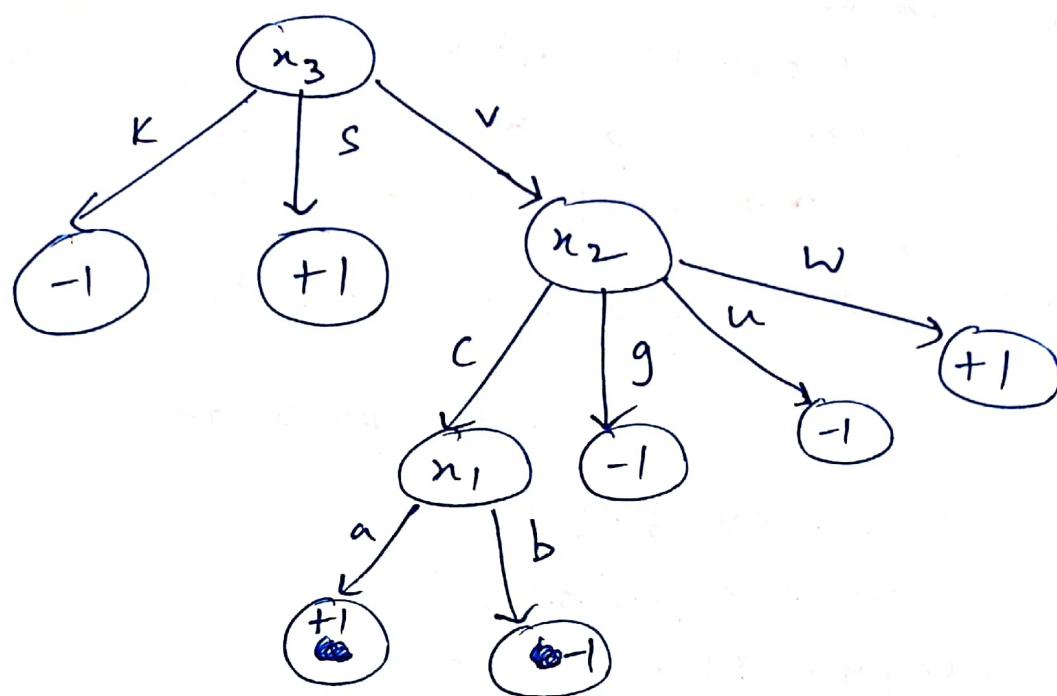
$$\begin{aligned} \text{Entropy}(S_3 | x_2) &= \frac{2}{5}(1) + \frac{1}{5}(0) + \frac{1}{5}(0) + \frac{1}{5}(0) \\ &= 0.4 \end{aligned} \quad \text{--- (4)}$$

$S_0$



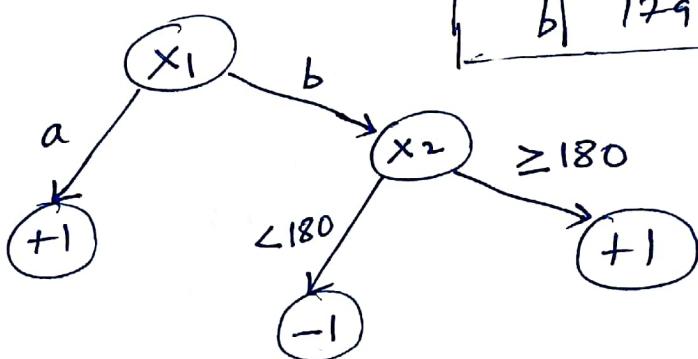
We split  $S_4$  with  $x_1$  which is only attribute left

Final Decision tree :



- (b).
- $x_1 \in \{a, b\}$
  - $x_2 \in \mathbb{N}$
  - $x_3 \in \{e, f\}$
  - $x_4 \in \{c, d\}$

i). Decision tree:



$y \in \{+1, -1\}$				
$x_1$	$x_2$	$x_3$	$x_4$	$y$
b	185	f	d	+1
b	180	f	c	+1
b	170	f	c	-1
b	140	e	d	-1
a	176	f	d	+1
b	129	f	d	-1

The decision tree has a 100% accuracy on the training data.

- (ii) .
- $x_1 = b, x_2 = 170, x_3 = f, x_4 = d \Rightarrow \text{pred} = -1 \text{ (correct)}$
  - $x_1 = a, x_2 = 150, x_3 = f, x_4 = d \Rightarrow \text{pred} = +1 \text{ (correct)}$
  - $x_1 = b, x_2 = 60, x_3 = f, x_4 = d \Rightarrow \text{pred} = -1 \text{ (incorrect)}$
- Test accuracy : 0.667 .