

clustering

Q. 1

$P_1(0.4, 0.53)$ ,  $P_2(0.22, 0.38)$ ,  $P_3(0.35, 0.32)$ ,  $P_4(0.26, 0.19)$   
 $P_5(0.08, 0.41)$ ,  $P_6(0.45, 0.30)$ .

A-

Use Single link (MIN) and Complete Link (MAX) technique to group the points into clusters.



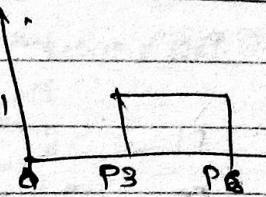
Proximity matrix based on Euclidean distance.

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$
$P_1$	0	0.22	0.22	0.37	0.34	0.23
$P_2$	0.24	0	0.15	0.20	0.14	0.25
$P_3$	0.22	0.15	0	0.15	0.28	0.11
$P_4$	0.37	0.20	0.15	0	0.29	0.22
$P_5$	0.34	0.4	0.28	0.29	0	0.39
$P_6$	0.23	0.25	0.11	0.22	0.39	0

Q1 →

A Using MIN (Single Link) technique.

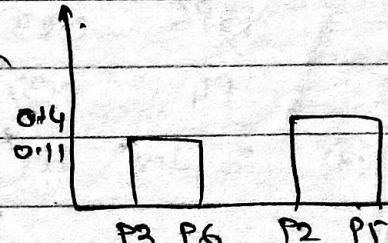
① As dist bet P3 & P6 is smallest, select it. & merge into first cluster.



New "dist" matrix -

	P1	P2	P3/P6	P4	P5
P1	0	0.24	0.22	0.37	0.34
P2	0.24	0	0.15	0.20	0.14
P3/P6	0.22	0.15	0	0.15	0.28
P4	0.37	0.20	0.15	0	0.39
P5	0.34	0.14	0.28	0.39	0

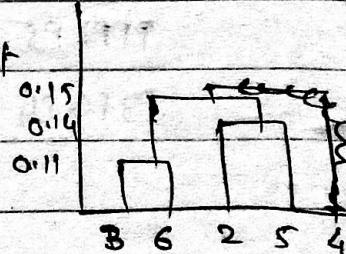
② As dist bet P2 & P5 is smallest, select it & merge into second cluster.



New "dist" matrix -

	P1	P2/P5	P3/P6	P4
P1	0	0.24	0.22	0.37
P2/P5	0.24	0	0.15	0.20
P3/P6	0.22	0.15	0	0.15
P4	0.37	0.20	0.15	0

③ As dist bet P2/P5 & P3/P6 is smallest, select them to merge into next cluster.

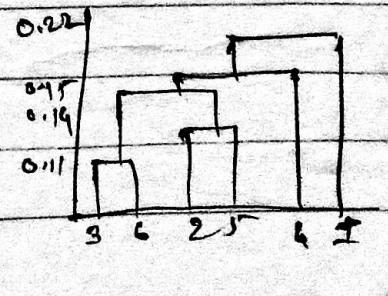


New "dist" matrix -

	P1	P2/P5	P3/P6	P4
P1	0	0.22	0.37	
P2/P5	0.22	0	0.15	0.20
P4	0.37	0.15	0	

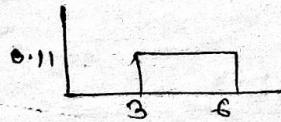
As dist bet {P2, P5, P3, P6} & {4} is smallest,

Select them to form next cluster.



	P1	P2/P5	P3/P6	P4
P1	0	22		
P2/P5	22	0		
P3/P6				
P4				

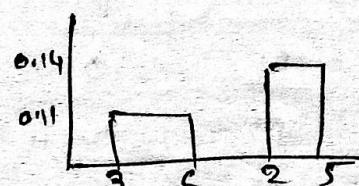
③ Using Complete Link (MAX) Technique  
 ① As dist' bet' 3 & 6 is smallest merge item



Dist' matrix.

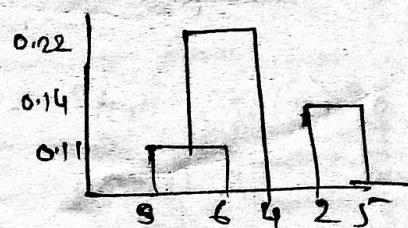
	P1	P2	P3P6	P4	P5
P1	0	0.24	0.23	0.37	0.34
P2	0.24	0	0.25	0.2	0.14
P3P6	0.23	0.25	0	0.22	0.39
P4	0.37	0.2	0.22	0	0.29
P5	0.34	0.14	0.39	0.29	0

② As dist' bet' P2 & P5 is smallest merge item



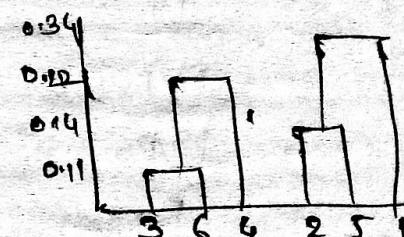
	P1	P2P5	P3P6	P4
P1	0	0.34	0.23	0.37
P2P5	0.34	0	0.39	0.29
P3P6	0.23	0.39	0	0.22
P4	0.37	0.29	0.22	0

③ As dist' bet' P3P6 & P4 is smallest merge item

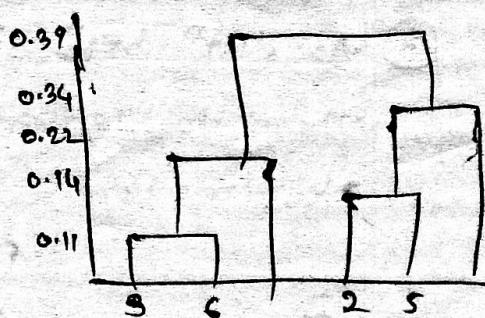


	P1	P2P5	P3P6P4
P1	0	0.34	0.37
P2P5	0.34	0	0.39
P3P6 P4	0.37	0.39	0

④ As dist' bet' P1 & P2P5 is smallest merge item



	P1 P2 P5	P3 P6 P4
P1 P2 P5	0	0.39
P3 P6 P4	0.39	0



CQ.2 : Use single link & complete link to cluster the data

(A)

	A	B	C	D	E
A	0	1	4	5	7
B	1	0	2	6	8
C	4	2	0	3	4
D	5	6	3	0	4
E	7	8	4	4	0

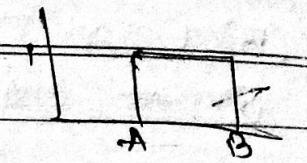
CQ.2(B) Use 2-means to cluster above data points  
using A, B as initial centroids.

C.2 (A) Using Single link technique -

(1) As dist bet A & B is smallest, merge them

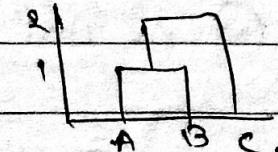
Dist matrix -

	AB	C	D	E
AB	0	2	5	7
C	2	0	3	4
D	5	3	0	4
E	7	4	4	0



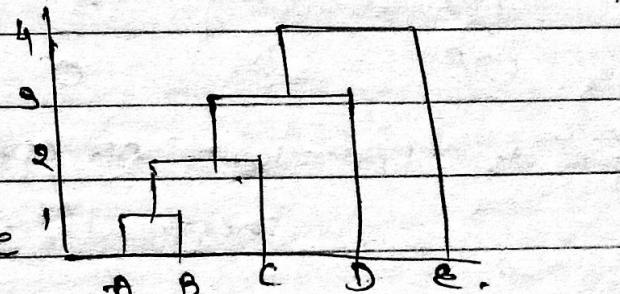
(2) As dist bet ABC & C is smallest merge them

	ABC	D	E
ABC	0	3	4
D	3	0	4
E	4	4	0



(3) As dist bet ABC & D is smallest merge them

	ABCD	E
ABCD	0	4
E	4	0



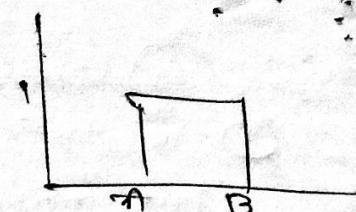
(4) Merge the final remaining node

with cluster formed.

(B) Using MAX (Complete link technique)

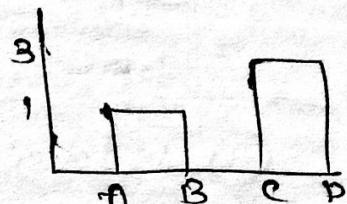
① As dist betn A & B is smallest, merge them.

	AB	C	D	E
AB	0	4	6	8
C	4	0	3	4
D	6	3	0	4
E	8	4	4	0



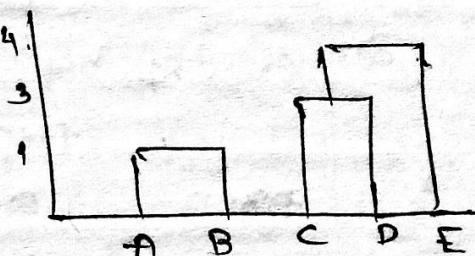
② As dist betn C & D is smallest merge them

	AB	CD	E
AB	0	6	8
CD	6	0	4
E	8	4	0

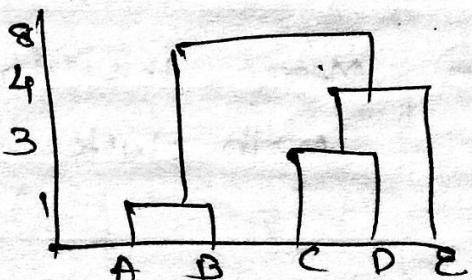


③ As dist betn E & CD is smallest merge them

	AB	CDE
AB	0	8
CDE	8	0



④ As dist betn AB, CDE is smallest merge them



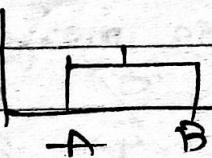
Q.2

Apply Average link algo for clustering

Ques. No.: 1 / 1  
पृष्ठ सं.:

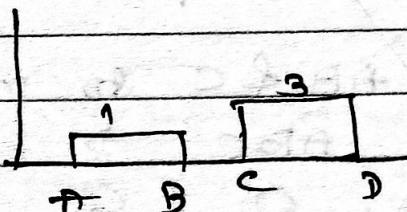
	A	B	C	D	E
A	0	1	4	5	7
B	1	0	2	6	8
C	4	2	0	3	4
D	5	6	3	0	4
E	7	8	4	4	0

Q.1 smallest dist is 1 b/w A & B so merge tree



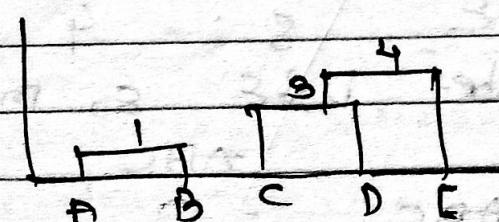
	AB	C	D	E
AB	0	3	5.5	7.5
C	3	0	3	4
D	5.5	3	0	4
E	7.5	4	4	0

Q.2 smallest dist is 3 b/w C & D so merge tree



	AB	CD	E
AB	0	4.25	7.5
CD	4.25	0	4
E	7.5	4	0

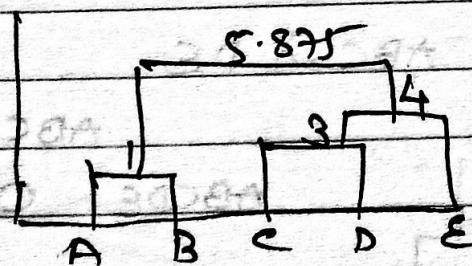
Q.3 smallest dist is 4 b/w CD & E so merge tree



	AB	CDE
AB	0	5.875
CDE	5.875	0

Q.4 last step so fuse AB & CDE together

ABCDE



cluster 1 → A, B

cluster 2 → C, D, E.

Use same data for Q4 & use K-means method to cluster the data points. Assume  $C_1(0.2, 0.3)$  &  $C_2(0.4, 0.4)$  as initial centroids.

(Q1, Q2) the data points: Assume  $C_1(0.2, 0.3)$  &  $C_2(0.4, 0.4)$  as initial centroids.

Dist' of points from  $C_1(0.2, 0.3)$ .

$$P_1(0.4, 0.53) \rightarrow d = \sqrt{(0.4 - 0.2)^2 + (0.53 - 0.3)^2} = 0.304$$

$$P_2(0.22, 0.38) \rightarrow d = \sqrt{(0.22 - 0.2)^2 + (0.38 - 0.3)^2} = 0.062$$

$$P_3(0.35, 0.32) \rightarrow d = \sqrt{(0.35 - 0.2)^2 + (0.32 - 0.3)^2} = 0.157$$

$$P_4(0.26, 0.19) \rightarrow d = \sqrt{(0.26 - 0.2)^2 + (0.19 - 0.3)^2} = 0.125$$

$$P_5(0.08, 0.41) \rightarrow d = \sqrt{(0.08 - 0.2)^2 + (0.41 - 0.3)^2} = 0.162$$

$$P_6(0.45, 0.3) \rightarrow d = \sqrt{(0.45 - 0.2)^2 + (0.30 - 0.3)^2} = 0.25$$

Dist' of points from  $C_2(0.4, 0.4)$ .

$$P_1(0.4, 0.53) \rightarrow d = \sqrt{(0.4 - 0.4)^2 + (0.53 - 0.4)^2} = 0.13$$

$$P_2(0.22, 0.38) \rightarrow d = \sqrt{(0.22 - 0.4)^2 + (0.38 - 0.4)^2} = 0.181$$

$$P_3(0.35, 0.32) \rightarrow d = \sqrt{(0.35 - 0.4)^2 + (0.32 - 0.4)^2} = 0.094$$

$$P_4(0.26, 0.19) \rightarrow d = \sqrt{(0.26 - 0.4)^2 + (0.19 - 0.4)^2} = 0.241$$

$$P_5(0.08, 0.41) \rightarrow d = \sqrt{(0.08 - 0.4)^2 + (0.41 - 0.4)^2} = 0.320$$

$$P_6(0.45, 0.3) \rightarrow d = \sqrt{(0.45 - 0.4)^2 + (0.3 - 0.4)^2} = 0.111$$

Points in cluster 1 with centroid  $C_1$  are  $P_2, P_4, P_5$ .

Points in cluster 2 with centroid  $C_2$  are  $P_1, P_3, P_6$ .

Need to check whether centroid of cluster changes or not.

$$\text{New centroid of cluster 1: } C_{1,\text{new}} = \frac{1}{3}(0.22 + 0.26 + 0.08) = 0.186$$

$$C_{1,\text{new},y} = \frac{1}{3}(0.38 + 0.19 + 0.41) = 0.326$$

$$C_{1,\text{new}} = (0.186, 0.326)$$

$$\text{New centroid for cluster 2: } C_{2,\text{new},x} = \frac{1}{3}(0.4 + 0.35 + 0.45) = 0.4$$

$$C_{2,\text{new},y} = \frac{1}{3}(0.53 + 0.32 + 0.3) = 0.383$$

$$C_{2,\text{new}} = (0.4, 0.303)$$

$$\text{SSF} = \sum_{i=1}^k \sum_{\alpha \in C_i} \text{dist}(C_i, \alpha)^2$$

where dist is standard Euclidean dist' b/w two objects.

$C_i$  is centroid of cluster i

$$\begin{aligned} &= (0.082)^2 + (0.125)^2 + (0.162)^2 + \\ &\quad (0.13)^2 + (0.094)^2 + (0.111)^2 = 0.08, \end{aligned}$$

Q.2 B  $k=2$ , so need to find out 2 clusters, using 3 means.

dist of points from centroid A  $\rightarrow$ .

$$d(A, B) \rightarrow \text{dist} =$$

$$\text{Kmean}(A, C) \rightarrow \text{dist} = 4$$

$$d(A, E) \rightarrow \text{dist} = 7$$

$$d(A, D) \rightarrow \text{dist} = 5$$

dist of points from centroid B  $\rightarrow$ .

$$d(B, C) = 2$$

$$d(B, D) = 6$$

$$d(B, E) = 8$$

from above distances, points in first cluster with centroid

A are  $\rightarrow A, E, D$ .

B are  $\rightarrow B, C$

New centroid for cluster 1 =

~~mean~~

$$1) V = 3 \leftarrow$$

$$2) V = 4 \leftarrow (2E)$$

$$3) V = 5 \leftarrow (n)$$

$$V = 3 \leftarrow$$

Apply Minimum Spanning Tree (MST) clustering technique to find out clusters in dataset.

पृष्ठ सं.: 1

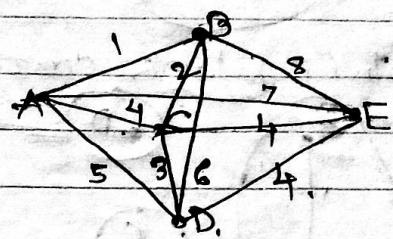
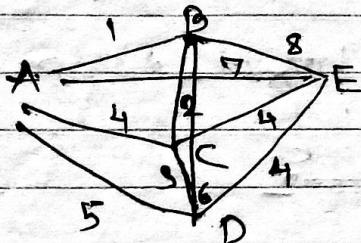
	A	B	C	D	E
A	0	1	4	5	7
B	1	0	2	6	8
C	4	2	0	3	4
D	5	6	3	0	4
E	7	8	4	4	0

Using Kruskal's algorithm find MST.

A MST of graph is subgraph that.

- (a) has no cycles i.e. is a tree.
- (b) contains all the nodes of graph
- (c) has minimum total edge weight of all possible trees.

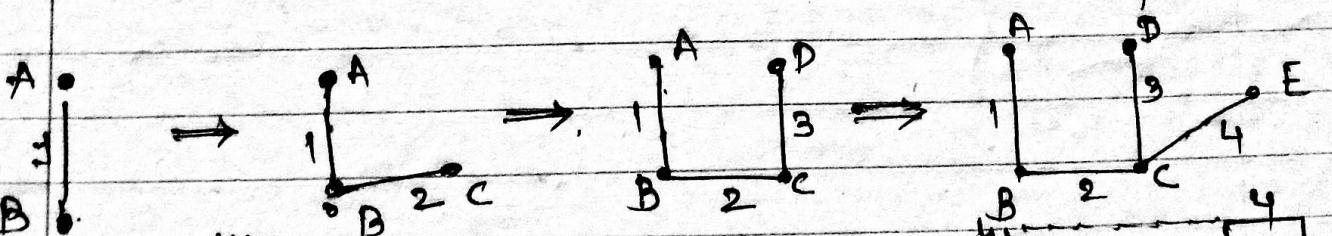
Step 1 → First construct graph for the points given.



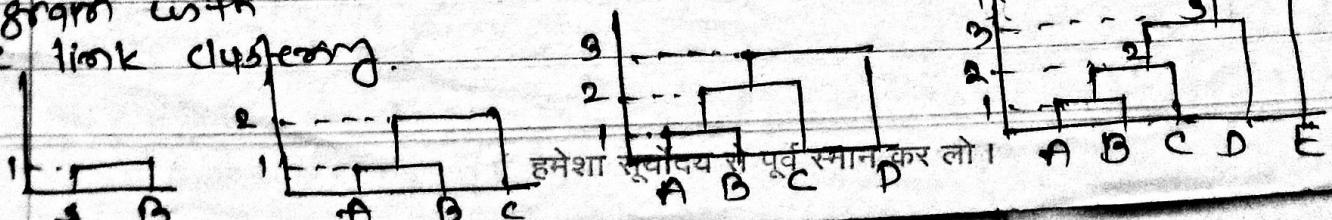
Step 2 → Arrange edges of graphs in increasing order of weight  
 $AB = 1 \quad BC = 2 \quad CD = 3 \quad CE, AC, DE = 4 \quad AD = 5 \quad BD = 6$   
 $AE = 7 \quad BE = 8$ .

Step 3 → Apply Kruskal's algorithm to build MST.

- (a) Sort all edges in increasing order of weight.
- (b) Pick smallest edge. Check if it forms cycle with S.T. formed so far. If cycle is not formed, include edge. Else discard it.
- (c) Repeat # b until there are  $V-1$  edges in tree, where  $V$  is number of vertices in graph.



Dendrogram with single link clustering.



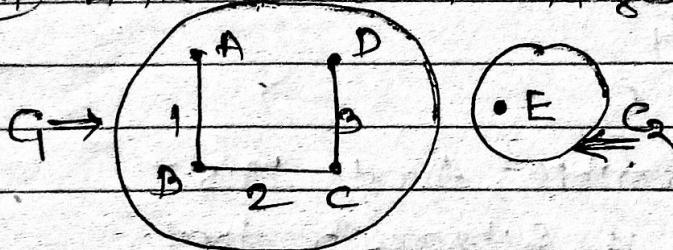
Repeat

Step 4

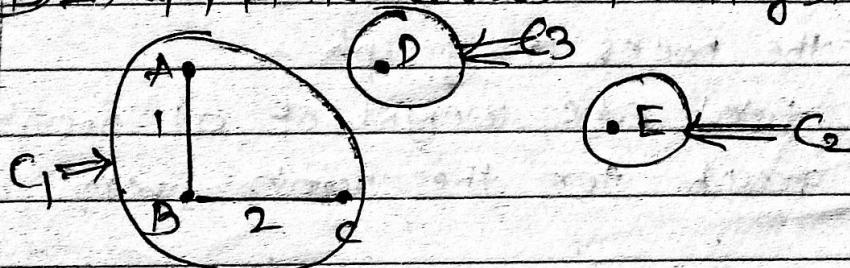
Create new cluster by breaking the link corresponding to the largest dissimilarity

Until only one singleton cluster remain

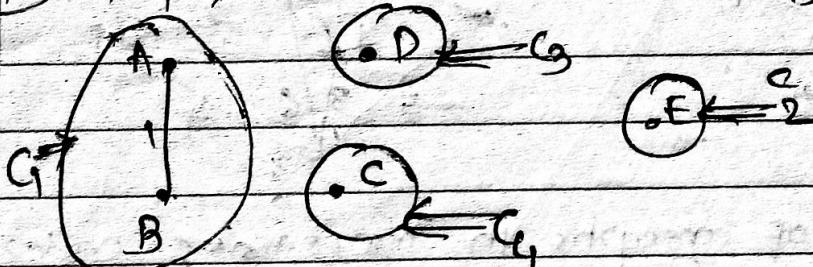
(4a) Link corner to largest dissimilarity  $CF = 4$



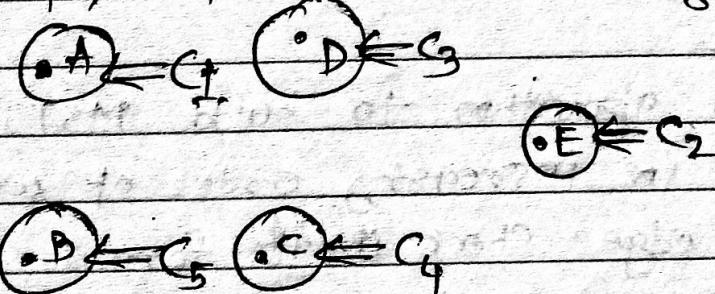
(4b) In  $C_1$ , Link corner to largest dissimilarity  $DC = 3$



(4c) In  $C_1$ , Link corner to largest dissimilarity  $BC = 2$



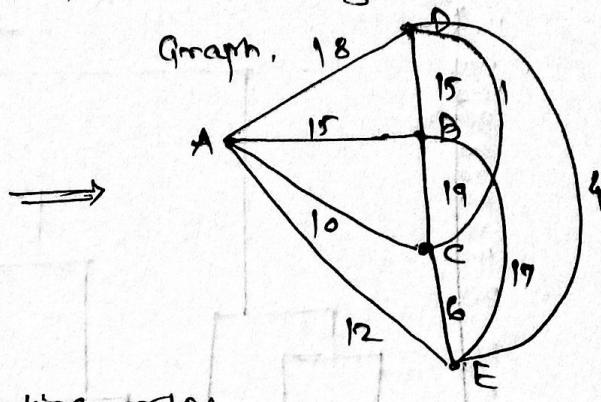
(4d) In  $C_1$ , Link corner to largest dissimilarity  $AB = 1$ .



(4e) All clusters are singletons, so no more splitting reqd

Q. Construct a graph showing all edges of following data, find MST for this graph. Is MST single link hierarchical clustering the same as that found using the tradition of single link algorithm?

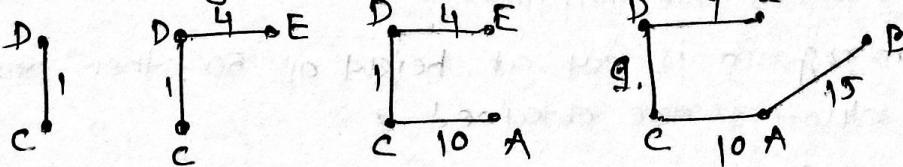
Item	A	B	C	D	E
A	0	15	10	18	12
B	15	0	19	15	17
C	10	19	0	1	6
D	18	15	1	0	4
E	12	17	6	4	0



Edges of graph in ascending order.

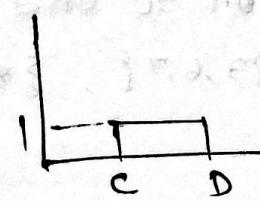
$$DE=4 \quad CE=6 \quad AC=10 \quad AE=12 \quad AB=15 \quad BE=17 \quad AD=18 \\ BC=19$$

MST using Kruskal's algorithm -



Traditional Single Link Algorithm -

	A	B	C	D	E
A	0	15	10	12	
B	15	0	15	17	
C	10	15	0	4	
E	12	17	4	0	



AB      CDE

A 0 15 10

B 15 0 15

CDE 10 15 0

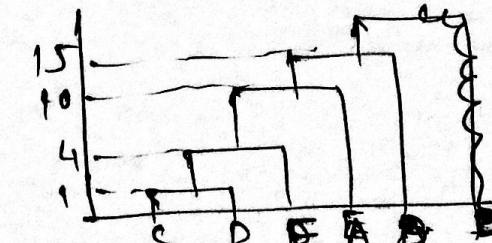
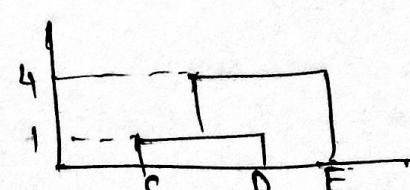
AB CDEA

A 0 15

B 15 0

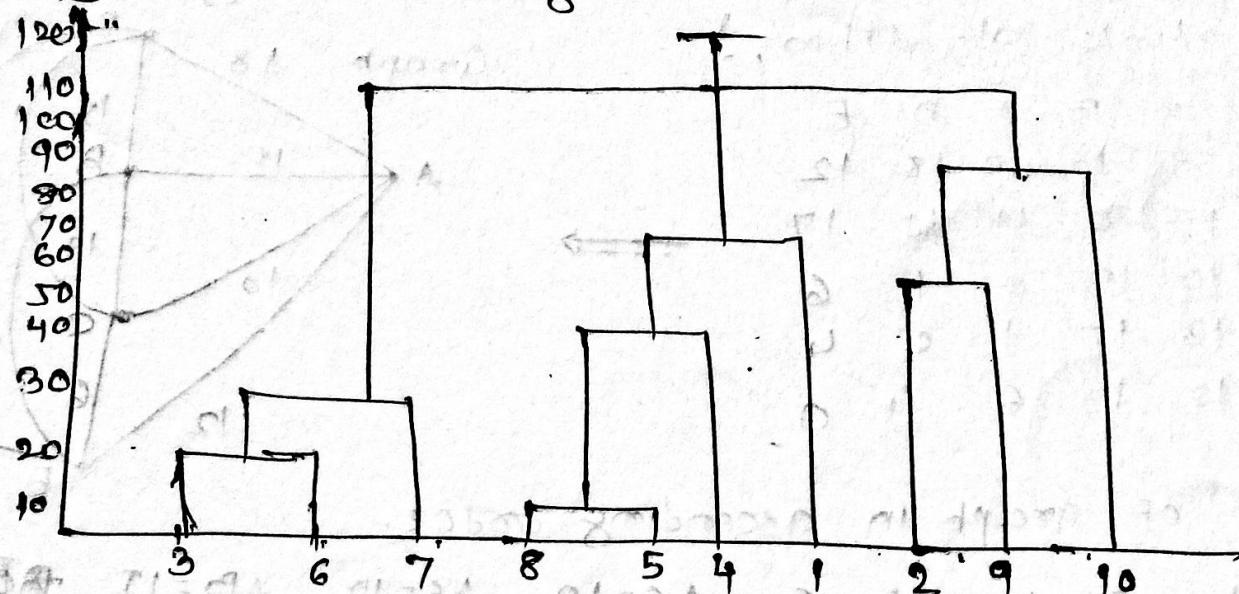
AB CDEA

AB CDEA



MST single link clustering produces same clusters as traditional single link clustering algorithm as the edges on MST are added in same order as clusters are merged.

Q. Consider the dendrogram obtained using hierarchical agglomerative clustering algorithm. Answer the following questions w.r.t. dendrogram.



- (a) How many natural clusters are there in dataset?
- (b) At what height the dendrogram must be cut, if three clusters are required?
- (c) If dendrogram is cut at height of 50, then how many clusters are obtained?

- 
- (a) 3     $\{3, 6, 7\}$ ,  $\{8, 5, 4, 1\}$  &  $\{2, 9, 10\}$ .
  - (b) b/w 90 and 110. excluding both values
  - (c) 6     $\{3, 6, 7\}$ ,  $\{8, 5, 4\}$ ,  $\{1\}$ ,  $\{2\}$ ,  $\{9\}$ ,  $\{10\}$