

700740862
Veera Venkata Satyanarayana
VXB08620

Assignment - 6

- Q.1 Find out clustering representations, & dendrograms using single, complete, and average link proximity function in Hierarchical clustering technique?

Point	X-Coordinate	Y-Coordinate
P ₁	0.0005	0.5306
P ₂	0.2148	0.3854
P ₃	0.3457	0.3156
P ₄	0.2652	0.1875
P ₅	0.0789	0.4139
P ₆	0.6548	0.3022

X-Y co-ordinate.

Table - 1

Distance Matrix.

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0.000	0.2357	0.2218	0.3688	0.3421	0.2347
P ₂	0.2357	0.000	0.1463	0.2042	0.1388	0.2540
P ₃	0.2218	0.1403	0.000	0.1513	0.2843	0.1100
P ₄	0.3688	0.2042	0.1513	0.000	0.2932	0.2216
P ₅	0.3421	0.1388	0.2843	0.2932	0.000	0.3941
P ₆	0.2347	0.2540	0.1100	0.2216	0.3941	0.000

Table - 2

By single link:

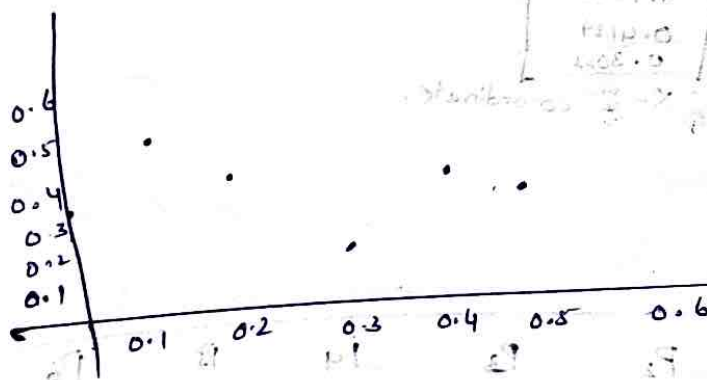
* For single link hierarchical clustering, the proximity of two clusters is minimum of the distance b/w any two points in 2 different clusters.

* The single link technique is good for non-elliptical shapes, but sensitive to noise & outliers.

* Applying single link techniques to our example data set of six points.

set of six 2 dimensional

points.



→ from table 1, we can observe distance b/w

P_3 & P_6 is 0.11.

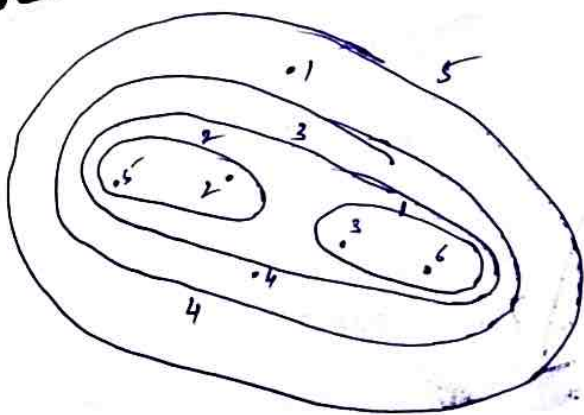
→ the height at which two clusters are merged can be represented as distance b/w two clusters.

distance b/w clusters $\{3, 6\}$ & $\{2, 5\}$ is given

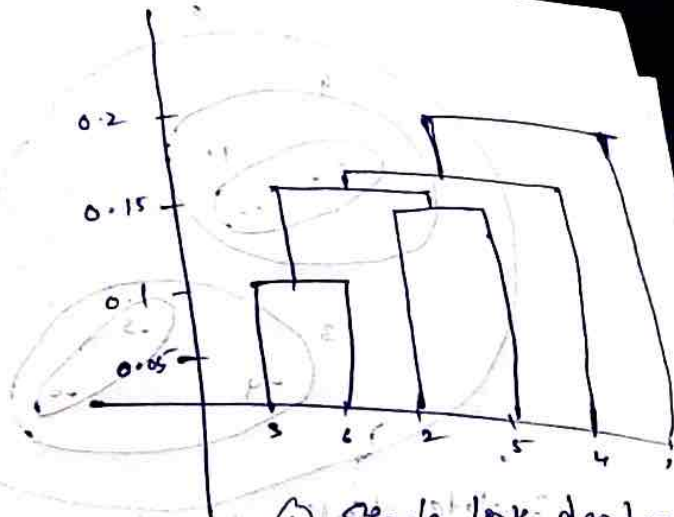
$$\text{by } \text{dist}(\{3, 6\}, \{2, 5\}) = \min(\text{dist}(3, 2), \text{dist}(6, 2), \text{dist}(3, 5), \text{dist}(6, 5))$$

$$\Rightarrow \min(0.15, 0.25, 0.28, 0.39)$$

$$\Rightarrow 0.15$$



Single link clustering



⑥ Single link dendrogram

Complete link.

→ In Complete link of hierarchical clustering, the proximity of two clusters is defined as the maximum of the distances b/w any two points in two different clusters.

→ Complete link is less susceptible to noise & outliers, but it can break large clusters & favours globular shapes.

→ Below fig shows results of Applying max to the sample data set of six points.

→ Here points 3 & 6 are merged first. {3, 6} is merged with {4} instead of {2, 5} or {1}. This is because

$$\text{dist}(\{3, 6\}, \{4\}) = \max(\text{dist}(3, 4), \text{dist}(6, 4))$$

$$= \max(0.15, 0.22)$$

$$= 0.22$$

$$\text{dist}(\{3, 6\}, \{2, 5\}) = \max(\text{dist}(3, 2), \text{dist}(6, 2), \text{dist}(3, 5), \text{dist}(6, 5))$$

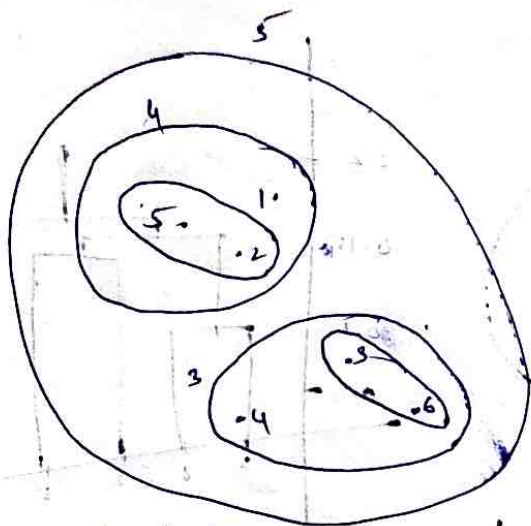
$$= \max(0.15, 0.25, 0.18, 0.39)$$

$$= 0.39$$

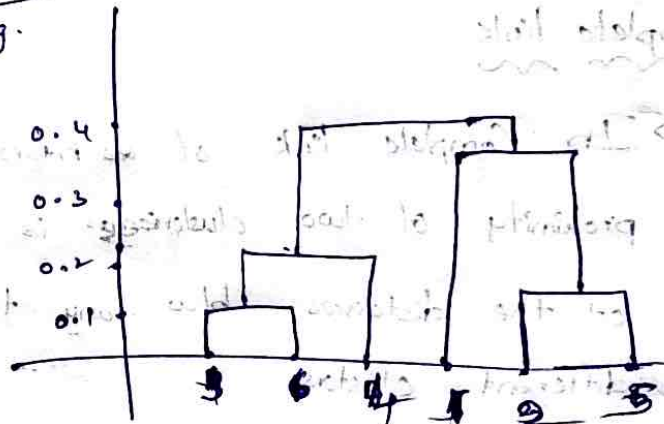
$$\text{dist}(\{3, 6\}, \{1\}) = \max(\text{dist}(3, 1), \text{dist}(6, 1))$$

$$= \max(0.22, 0.25)$$

$$= 0.25$$



Complete link clustering.



Complete link dendrogram.

Average link

Below figure shows results after applying the group Average approach to a sample data to six cluster points.

→ we calculate the distance b/w some clusters.

$$\rightarrow \text{proximity} \Rightarrow \text{proximity}(C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} \text{proximity}(x, y)}{m_i \times m_j}$$

$$\text{dist}(\{3, 6, 4\}, \{1\}) = (0.22 + 0.37 + 0.23) / (3 \times 1) = 0.28$$

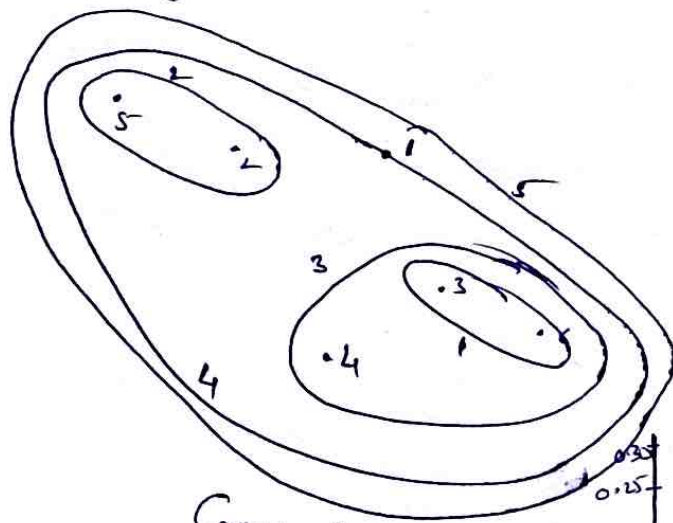
$$\text{dist}(\{2, 5\}, \{1\}) = (0.24 + 0.34) / (2 \times 1) = 0.29$$

$$\text{dist}(\{3, 6, 4\}, \{2, 5\}) = (0.15 + 0.28 + 0.25 + 0.39) / (3 \times 2) = 0.26$$

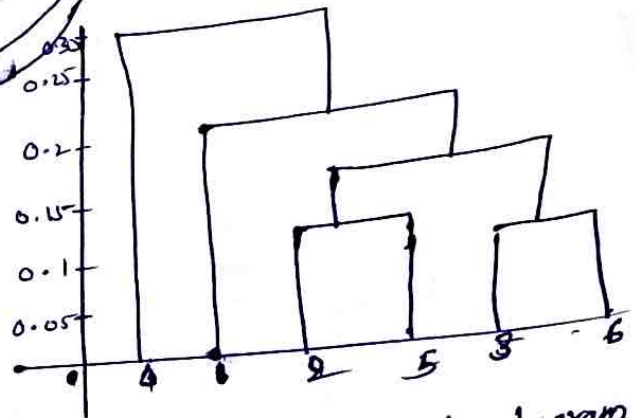
$$\text{dist}(\{1, 2\}, \{1, 3\}) = (0.15 + 0.20 + 0.20) / (3 \times 2) = 0.26$$

$$\text{dist}(\{1, 2\}, \{1, 3\}) = (0.15 + 0.20 + 0.20) / (3 \times 2) = 0.26$$

Here Because $\text{dist}(\{3, 6, 4\}, \{2, 5\})$ is smaller than $\text{dist}(\{3, 6, 4\}, \{1\})$ and $\text{dist}(\{2, 5\}, \{1\})$ clusters $\{3, 6, 4\}$ and $\{2, 5\}$ are merged at the fourth stage.



Group Average clustering.



→ Average version of hierarchical clustering, the proximity of two clusters is defined as the average pairwise proximity among all pairs of points in the different clusters.

proximity $\text{proximity}(C_i, C_j)$ of clusters C_i and C_j which are of size m_i and m_j respectively is

$$\text{proximity}(C_i, C_j) = \frac{\sum_{x \in C_i} \text{proximity}(x, y)}{m_i \times m_j}$$