

Question-1

Find out clustering representations & Dendrogram using single, complete and Average link proximity function in Hierarchical clustering technique.

point	x-coordinate	y co-ordinate
P1	0.04005	0.5306
P2	0.2148	0.3854
P3	0.3457	0.3156
P4	0.2652	0.1875
P5	0.0789	0.4139
P6	0.4548	0.3022

Table 1

X-Y coordinates.

Distance matrix:

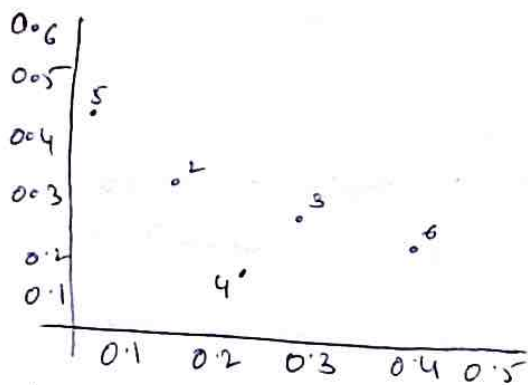
	P1	P2	P3	P4	P5	P6
P1	0.00	0.2357	0.2418	0.3688	0.3421	0.2347
P2	0.2357	0.000	0.1483	0.2042	0.1388	0.2540
P3	0.2418	0.1483	0.000	0.1513	0.2843	0.1100
P4	0.3688	0.2042	0.1513	0.000	0.2932	0.2216
P5	0.3421	0.1388	0.2843	0.2932	0.000	0.3921
P6	0.2347	0.2540	0.1100	0.2216	0.3921	0.000

Table 2

By Single link:-

- For Single link hierarchical clustering, the proximity of two clusters is minimum of the distance between 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 technique to our example data set of Simpson's set of six 2 dimensional points.

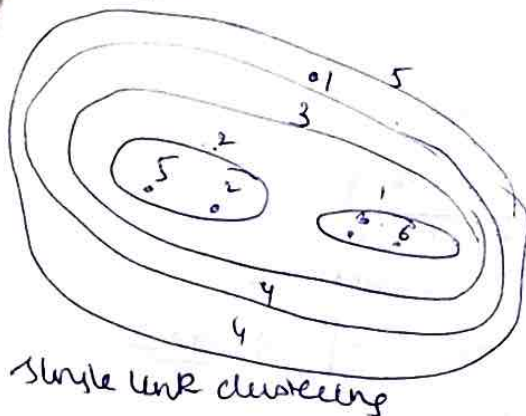


→ From table 1, we can observe distance between P_2 & P_6 is 0.11

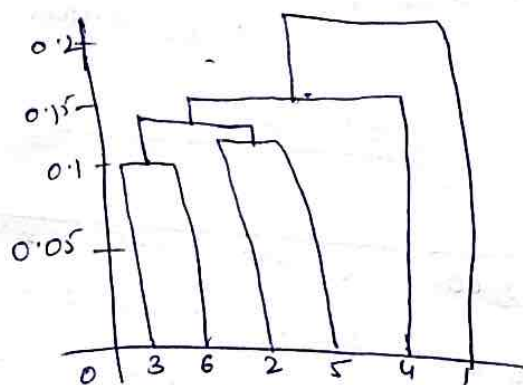
- The height at which two clusters are merged can be represented as distance between two clusters.

distance between cluster $\{3, 6\}$ & $\{2, 5\}$ is given 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



(b) single link dendrogram

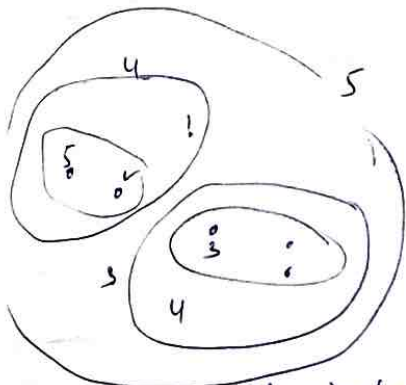
Complete link

- In complete link or hierarchical clustering, the proximity of two clusters is defined as the maximum of the distance between any two points in two different clusters.
- Complete link is less susceptible to noise & outliers, but it can break large clusters & it favours globular shapes.
- Below is shown result of applying max to the sample data set of six points.
- Here points 3 and 6 are merged first.
- $\{3, 6\}$ is merged with $\{5\}$ instead of $\{2, 5\}$ or $\{1\}$ then because

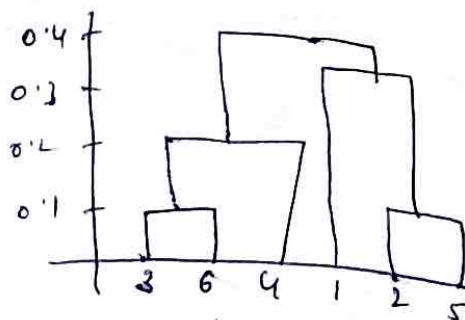
$$\text{dist}(\{3,6\}, \{4\}) = \max(\text{dist}(3,4), \text{dist}(6,4)) \\ = \max(0.15, 0.21) \\ = 0.22$$

$$\text{dist}(\{3,6\}, \{1,5\}) = \max(\text{dist}(3,1), \text{dist}(6,1), \text{dist}(3,5), \text{dist}(6,5)) \\ = \max(0.15 + 0.25, 0.28, 0.39) \\ = 0.39$$

$$\text{dist}(\{3,6\}, \{1\}) = \max(\text{dist}(3,1), \text{dist}(6,1)) \\ = \max(0.22, 0.23) = 0.23$$



complete link dendrogram



complete link dendrogram

Average Link:-

Below figure shows results after applying the group mean approach to sample data of six points.

The calculation of the distance between some clusters.

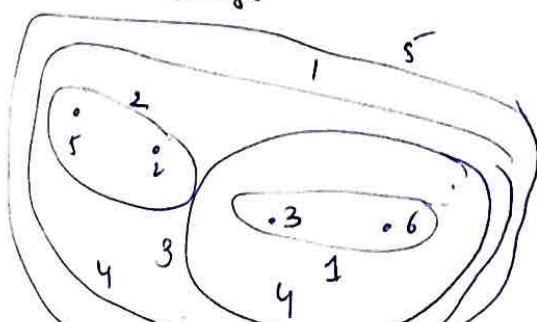
$$\text{proximity} \Rightarrow \text{proximity}(C_i, C_j) = \frac{\sum_{i \in C_i} \sum_{j \in C_j} \text{proximity}(i, j)}{m_i \times m_j}$$

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

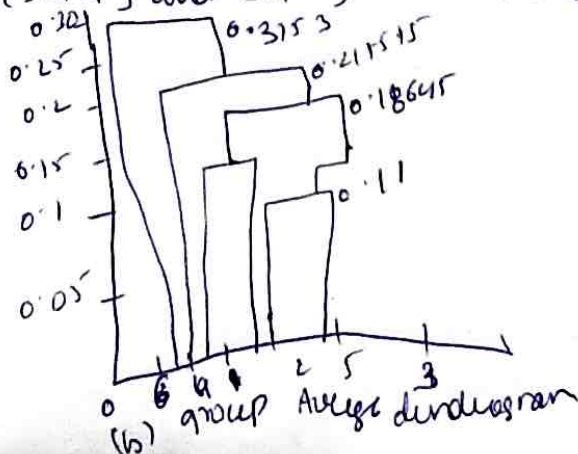
$$\text{dist}(\{1,5\}, \{1\}) = (0.14 + 0.34) / (2 \times 1) = 0.24$$

$$\text{dist}(\{3,6\}, \{1,5\}) = (0.15 + 0.28 + 0.25 + 0.39 + 0.20 + 0.29) / (3 \times 2) \\ = 0.26$$

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



group average dendrogram



(b) group average dendrogram

→ Another 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 proximity (C_i, C_j) of clusters C_i and C_j which are of size m_i and m_j respectively is

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