

### Assignment-3

Q1

It is not a valid distance as it does not satisfy triangle inequality.

$$d(x, y) = |x - y|^2$$

$$|x - z| \leq |x - y| + |y - z|$$

Squaring  $|x - z|^2 \leq |x - y|^2 + |y - z|^2 + 2|x - y||y - z|$

This is not guaranteed.  $\Rightarrow$  extra term

Q2

$$d_{max}^{PS} = \max_{y \in C} (x, y) = d(x, x_8)$$

$$(6, 4) \in (1.5, 1.5), d_{max} = \sqrt{(6 - 1.5)^2 + (4 - 1.5)^2} = 5.15$$

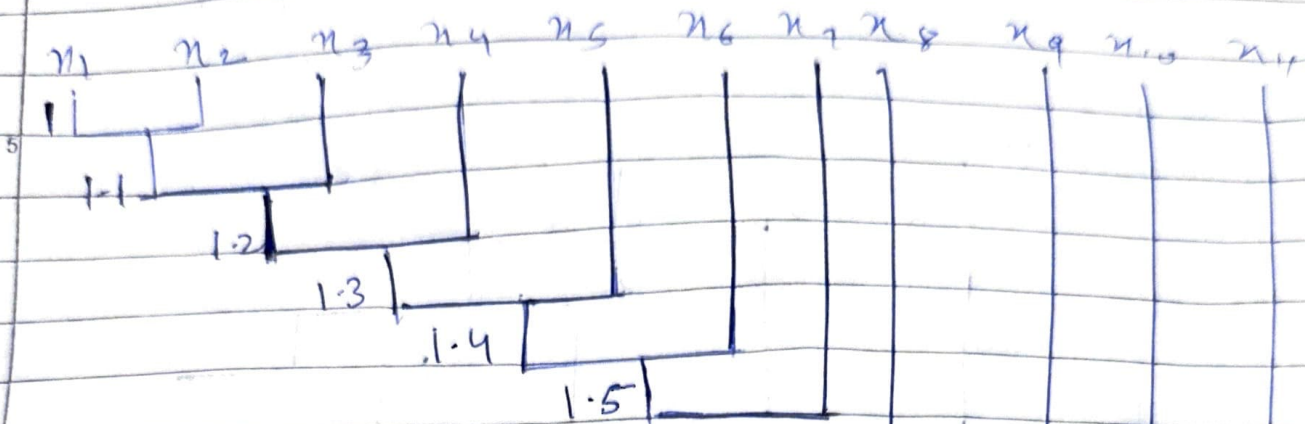
$$d_{min}^{PS}(x, C) = \min_{y \in C} (x, y) = d(x, x_3)$$

$$(6, 4) \in (3.5, 3), d_{min} = \sqrt{(6 - 3.5)^2 + (4 - 3)^2} = 2.69$$

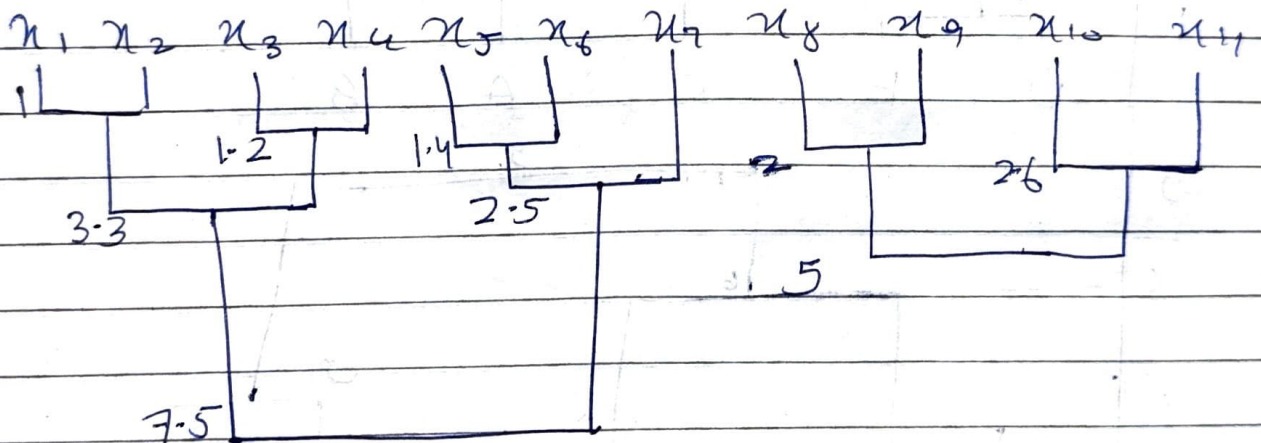
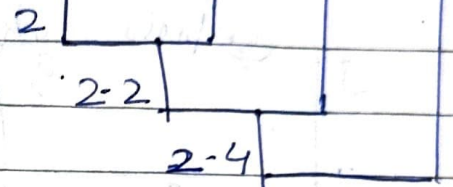
$$d_{avg}^{PS}(x, C) = \frac{1}{n_C} \sum_{y \in C} d(x, y)$$

$$= \frac{1}{8} \sum_{i=1}^8 d(x, x_i) = 4.33$$

Q3



Dissimilarity  
dendrogram based  
or Single link algorithm

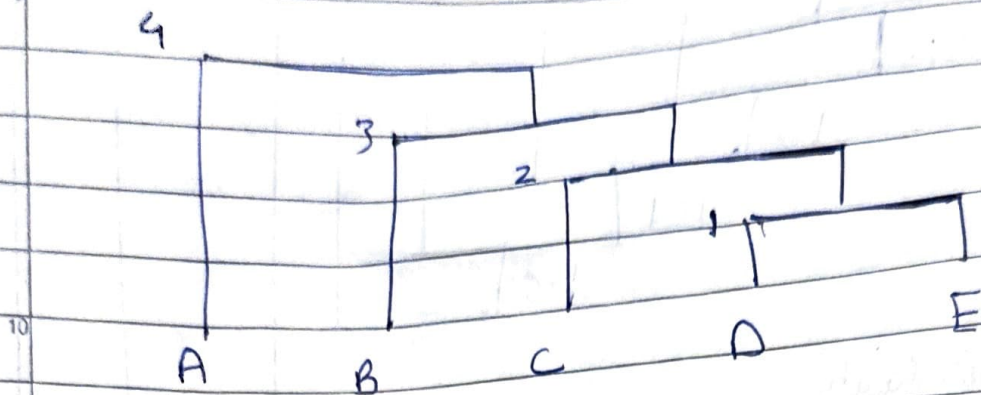


Complete link algorithm

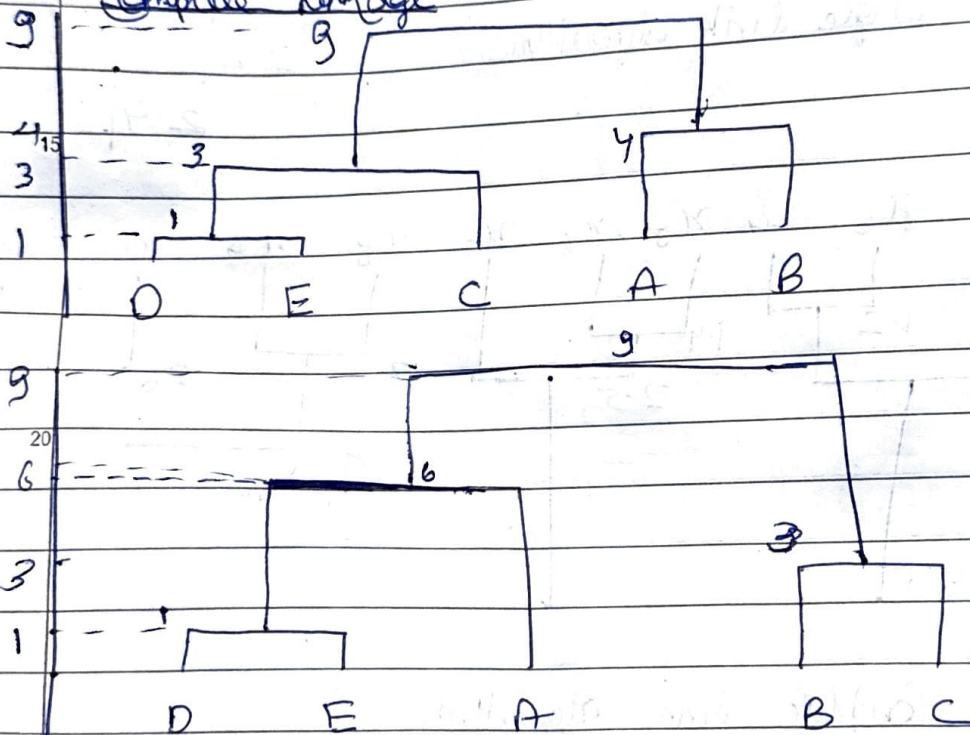


Q-4

# Single Linkage Dendrogram



## Complete Linkage



Q5  $C_1 = (2, 3)$   $C_2 = (5, 8)$   $C_3 = (9, 4)$

$(1, 2) \rightarrow C_1$

(A)  $(3, 4) \rightarrow C_1$

$(6, 7) \rightarrow C_2$

$(8, 3) \rightarrow C_3$

$(5, 5) \rightarrow C_2$

$C_1$  - Mean  $C'_1 = (1+3, 2+4) = (2, 3)$

$C_2$  - Mean  $C'_2 = (6+5, 7+5) = (5.5, 6)$

$C_3$  - Mean  $C'_3 = (8, 3)$

(b) Distortion -  
Initial -

$$\begin{array}{l|l} \text{Op} \text{ @ } \text{ @ } & 2 \text{ from } (1, 2) + 2 \text{ from } (3, 4) \\ \text{ @ } \text{ @ } & + 2 \text{ from } (6, 7) + 2 \text{ from } (8, 3) \\ & + 1 \text{ from } (5, 5) = 17 \end{array}$$

Final -

$$2 \text{ from } (1, 2) + 2 \text{ from } (3, 4) + 1.25 \text{ from } (6, 7) + 0 \text{ from } (8, 3) + 1.25 \text{ from } (5, 5) = 6.5$$

distortion decreases.



Q-6

Gaussian density -

$$\ln N(\mathbf{y} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$= -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln|\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{y}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{y}_n - \boldsymbol{\mu}_k)$$

Maximizing with respect to  $\boldsymbol{\Sigma}_k$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln(N) = -\frac{D}{2} \boldsymbol{\Sigma}_k^{-1} - \frac{1}{2} \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{y}_n - \boldsymbol{\mu}_k) (\mathbf{y}_n - \boldsymbol{\mu}_k)^T$$

$$\Rightarrow \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{y}_n - \boldsymbol{\mu}_k) (\mathbf{y}_n - \boldsymbol{\mu}_k)^T$$

Maximizing with respect to  $\pi_k$

we use -  $\sum_{n=1}^N \sum_{k=1}^K \ln(\pi_k) \gamma(z_{nk})$

under constraint of  $\sum_{k=1}^K \pi_k = 1$

So,  $\sum_{n=1}^N \gamma(z_{nk}) \ln(\pi_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$

diff - ,  $= \frac{N_k}{\pi_k} + \lambda = 0$

$$\Rightarrow \pi_k = -\frac{N_k}{\lambda}$$

$$\sum_{k=1}^K \pi_k = \sum_{k=1}^K -\frac{N_k}{\lambda} \Rightarrow \pi_k = \frac{N_k}{N}$$

Q-7 Conditional density

$$p(x_k | x_a) = \frac{p(x_a, x_k)}{p(x_a)}$$

$$= \sum_{k=1}^K \left( \frac{\pi_k p(x_k | x_a)}{\sum_{j=1}^K \pi_j p(x_j | x_a)} \right) p(x_a | x_k, R)$$

Mixing coefficient - 
$$\frac{\pi_k p(x_k | x_a)}{\sum_{j=1}^K \pi_j p(x_j | x_a)}$$

Component densities -  $p(x_k | x_a, R)$

Q-8 (a) Complete log-likelihood

$$\log p(x, z | \theta) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left[ \log \pi_k - \frac{1}{2} \log | \Sigma_k | \right] - \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)$$

(b) MLE -

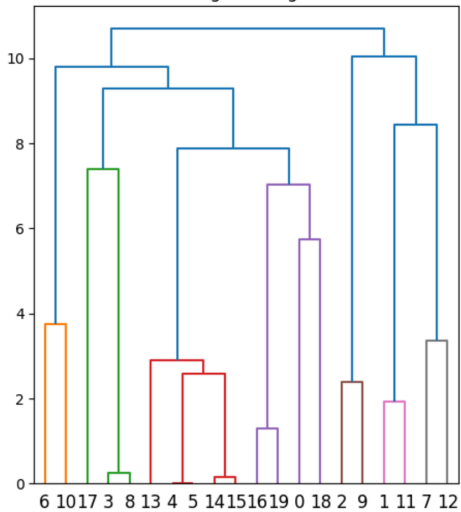
Maximize  $\sum_k \pi_k \log \pi_k$  under  $\sum_k \pi_k = 1$   
Using Lagrange multiplier -

$$\pi_k = \frac{n_k}{N}, \quad \pi_k = \frac{1}{N} \sum_{n=1}^N z_{nk}$$

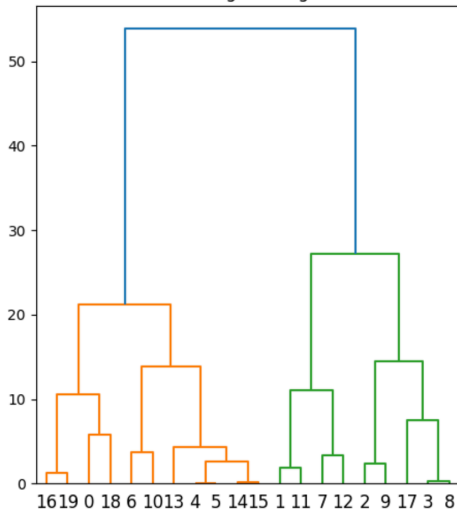
Covariance  $\Sigma_k$

$$\Sigma_k = \frac{1}{n_k} \sum_{n=1}^N z_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

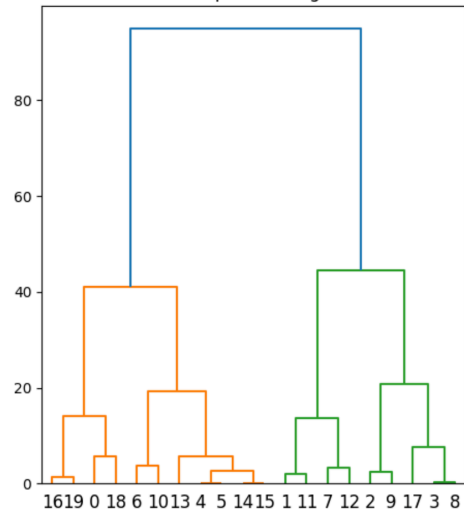
Single Linkage



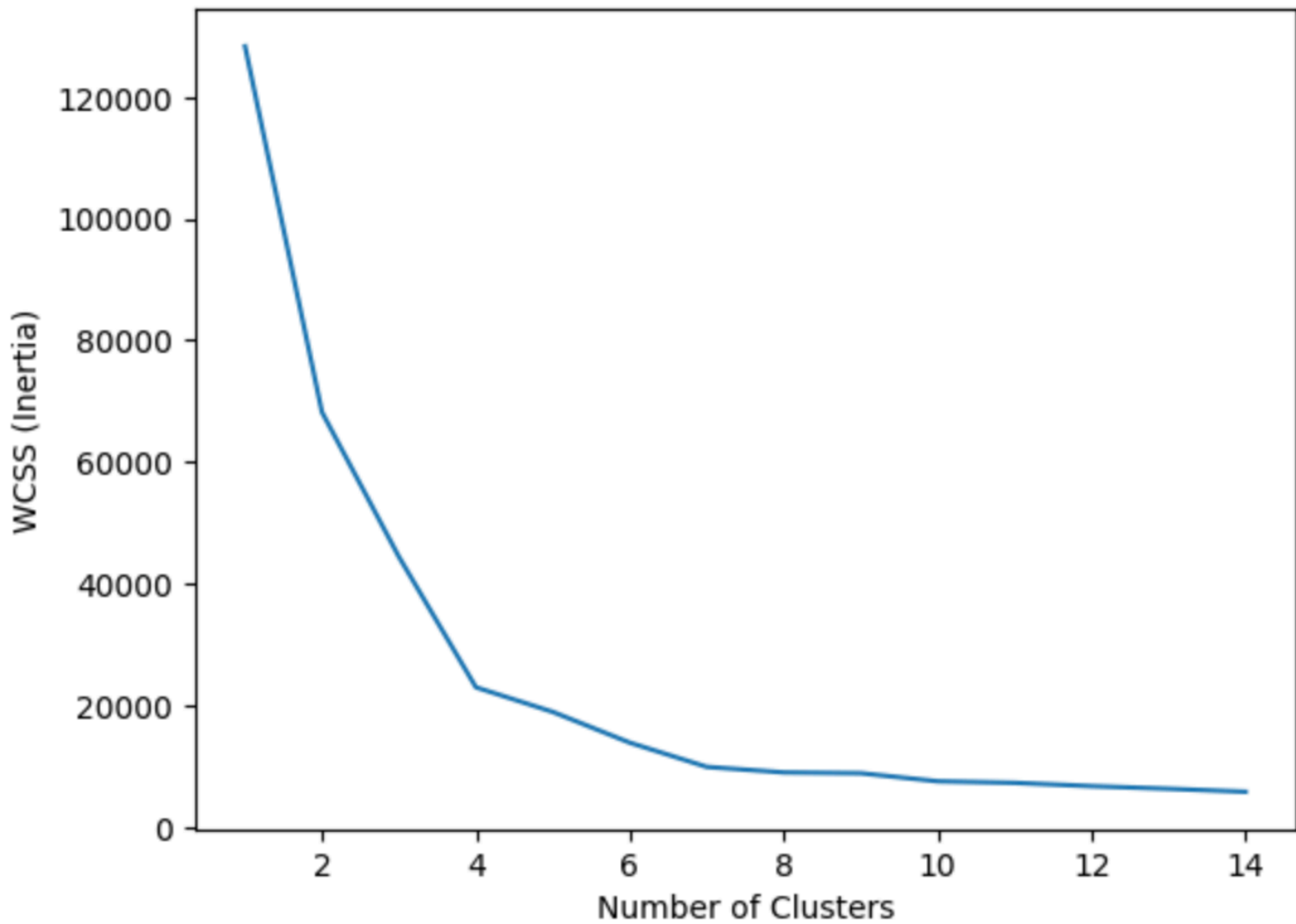
Average Linkage



Complete Linkage



# Elbow Method For Clusters





Cluster Plot

