

# \* Clustering

$$X_{n \times d} \rightarrow X = \bigcup_{j=1}^k X^{(j)} \xrightarrow{\text{Union of Clusters}} \text{Union of Clusters}$$

$k \rightarrow \text{clusters} \rightarrow 1 \text{ to } k$

## k-means algorithm

Find  $k$  disjoint clusters  $X^{(1)}, \dots, X^{(k)}$  such that, if  $\mu_1, \dots, \mu_j$  is the average point in cluster  $X^{(j)}$ , the following quantity is minimized

Minimize Distortion

$$\text{distortion} = \sum_{j=1}^k \sum_{x \in X^{(j)}} \|x - \mu_j\|^2$$

$\rightarrow \text{arg in } j^{\text{th}} \text{ cluster}$   
 $\rightarrow \text{mean of points in each cluster}$   
 $\swarrow$   
distance (euclidean)

## Basic Algorithm Steps

1/ Random Initialization: Find  $k$  random centers

$\mu_1, \dots, \mu_k \rightarrow \text{avg of clusters}$

Repeat

2/ Assignment Step: Assign each point to the nearest center  
That makes  $k$  clusters

$X_1, \dots, X_k$

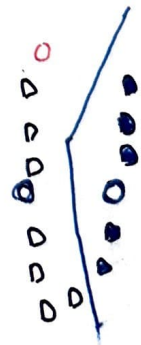
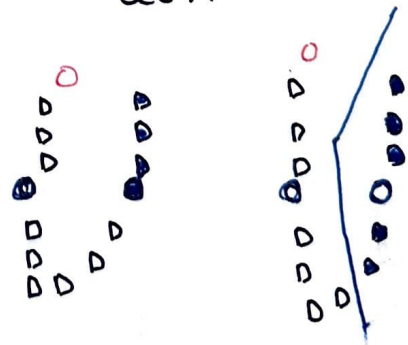
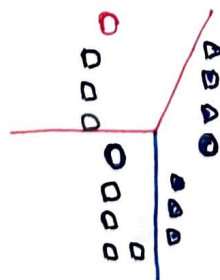
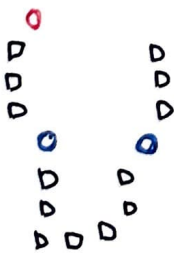
3/ Update Step: Replace the previous centers with the average point in each cluster

find new centers

Repeat this until no change in cluster memberships

$$\mu_j = \frac{1}{|X^{(j)}|} \sum_{x \in X^{(j)}} x$$

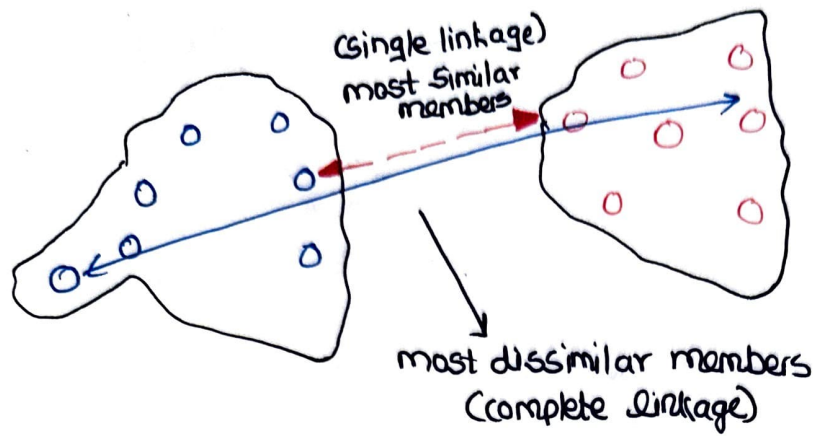
Voronoi diagram of space



## Issues

- The distortion objective is appropriate for spherical and coherent clusters
  - The algorithm is not guaranteed to find the best clusters for the given objective
  - The algorithm can be very slow to converge, so in practice we have to implement stopping criteria
- multiple initialization  
 $2\sqrt{n}$  steps

# \* Agglomerative Hierarchical Clustering



## Distance between Clusters

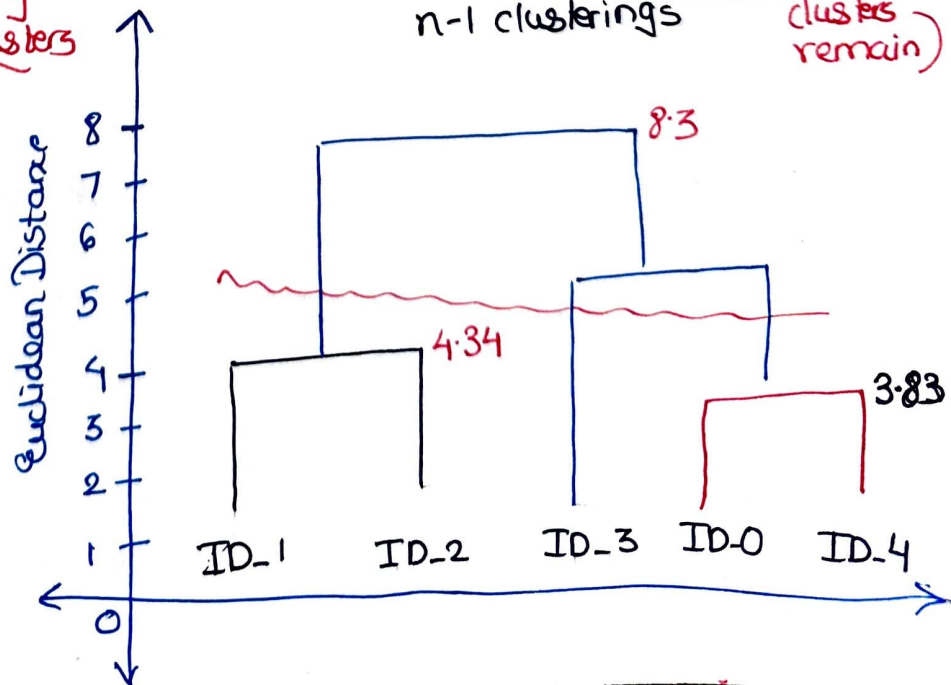
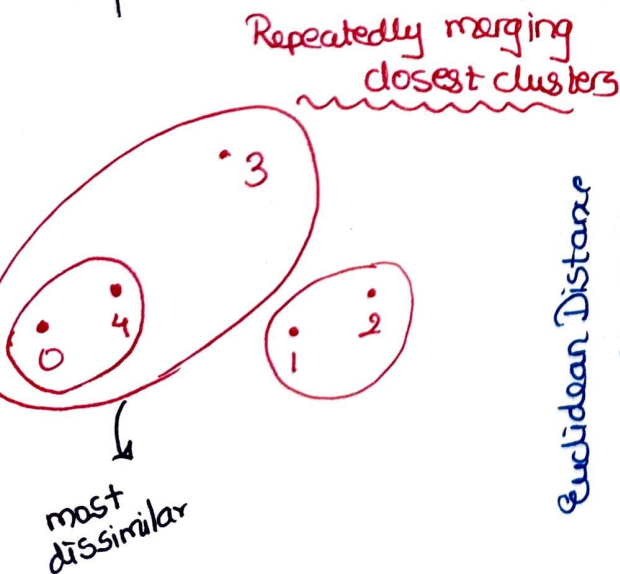
2 options → one is most similar  
other is most dissimilar

	ID-0	ID-1	ID-2	ID-3	ID-4
ID-0	0.00	4.97	5.51	5.89	3.83
ID-1	4.97	0.00	4.34	5.10	6.69
ID-2	5.51	4.34	0.00	7.24	8.31
ID-3	5.89	5.10	7.24	0.00	4.38
ID-4	3.83	6.69	8.31	4.38	0.00

same obvious

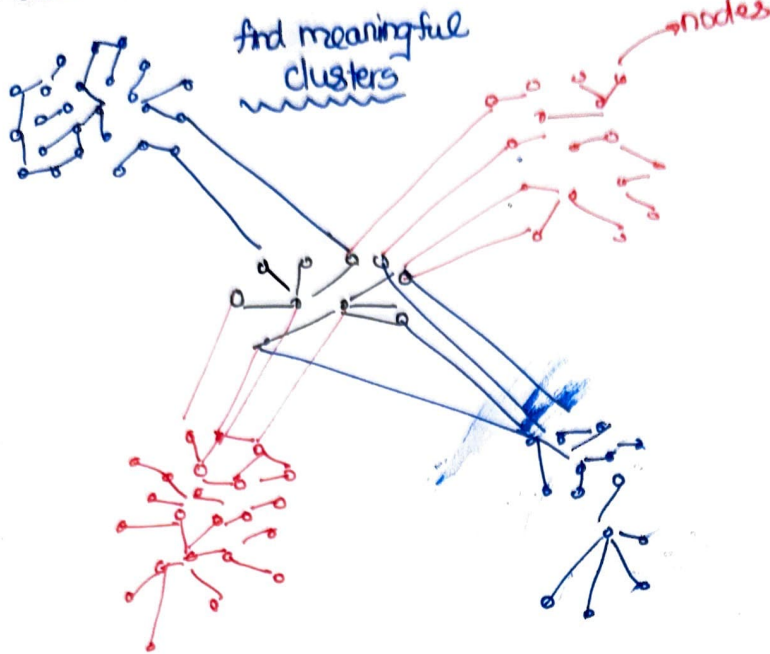
- 1) Start with individual clusters
- 2) Compute pairwise distances (Single, complete linkage)
- 3) Merge closest clusters
- 4) Update the distances (between newly formed clusters)
- 5) Repeat (until one/ desired no of clusters remain)

n-1 clusterings



# \* How to Cluster Networks :-

Clustering Networks → Social media, Genomes. etc

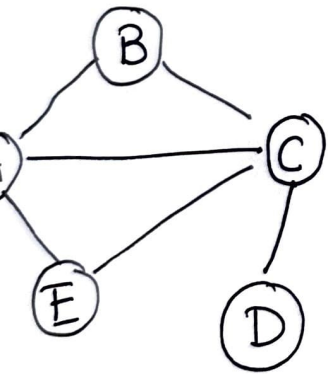


$$g(C) = \frac{\text{\# connections to other clusters}}{\text{\# connections inside cluster}}$$

fee objective ↗

Fi Value  
find clusters with small  
Fi value

Undirected  
Graph



Adjacency  
matrix

X	A	B	C	D	E
A	0	1	1	0	1
B	1	0	1	0	0
C	1	1	0	1	1
D	0	0	1	0	0
E	1	0	1	0	0

degrees

$$D = \begin{pmatrix} 3 \\ 2 \\ 4 \\ 1 \\ 2 \end{pmatrix}$$

diagonal matrix.

$$M = D^{-1}A$$

k eigenvalues and eigenvectors.

$\lambda_1, v_1$

k top eigenvalues } put into columns.

$(\lambda_i, v_i)$

$$Y = \begin{vmatrix} v_1 & v_2 & \dots & v_n \end{vmatrix}$$

matrix of points