

# Clustering

↳ grouping similar data points together into clusters based on characteristics/features

↳ unsupervised (no target variable)

↳ Goal:-

- identify natural grouping in data
  - minimize distance/dissimilarity between data points within the same cluster
  - maximize distance/dissimilarity between data points in different clusters
- group smaller distance together, larger distance apart

## Measures for Distance

↳ Euclidean distance

- common measurement
- $\sqrt{(x-x_i)^2 + (y-y_i)^2}$

↳ Manhattan distance

- absolute difference
- image/pattern recognition
- $|x_i - x_j| + |y_i - y_j|$

↳ Cosine similarity

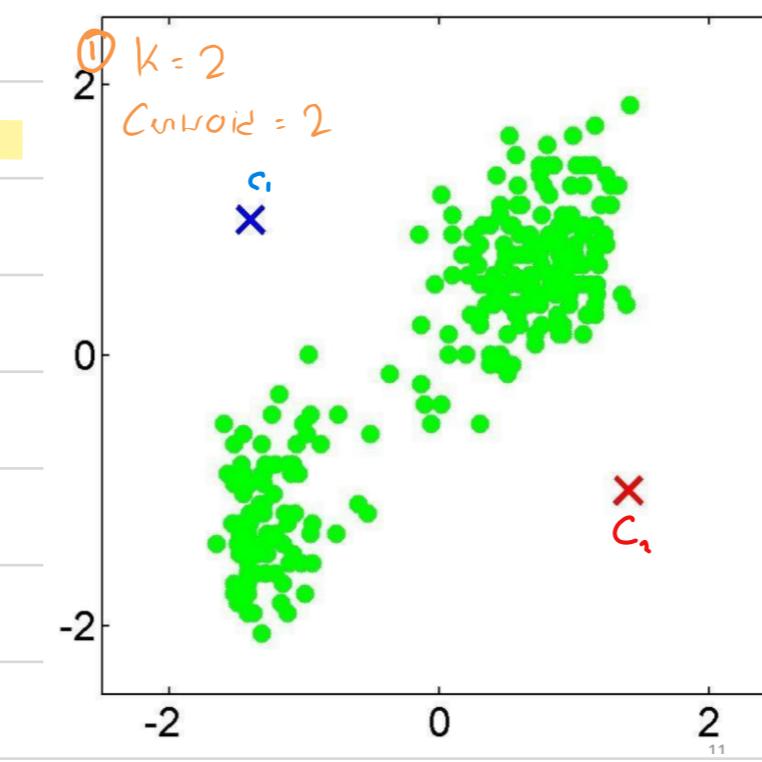
- based on cosine angle
- language recommendation
- $\cos \theta$

## K-Means

↳ clustering algorithm that partitions dataset into  $k$ -distinct clusters, where  $k$  is specified

↳ 1. Initialization:

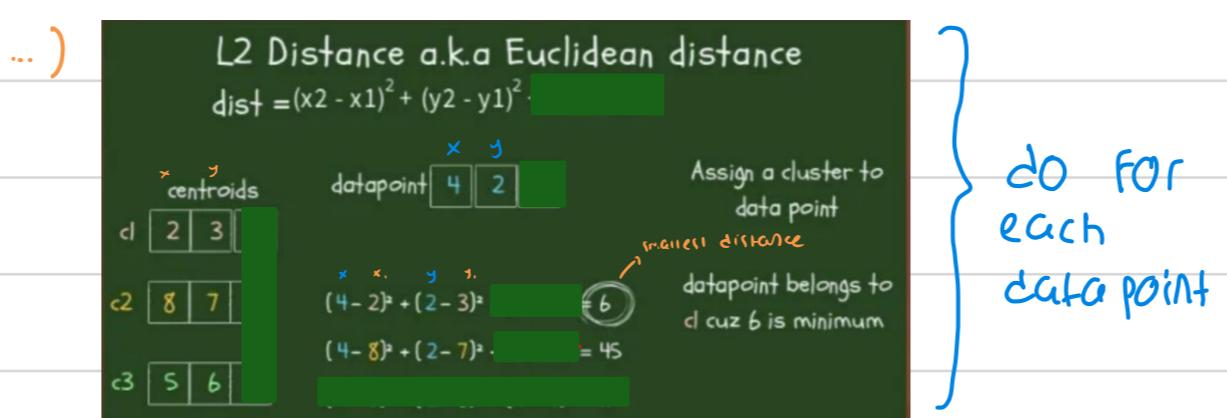
- ↳ select  $k$  random points as centroid
- by random selection
- ↳ randomly select from dataset
- by sort & split
- ↳ sort dataset & split into  $k$  portion



↳ 2. Assignment:

↳ Assign each data point to the closest centroid cluster, by distance measure (Euclidean, Manhattan...)

- Euclidean distance

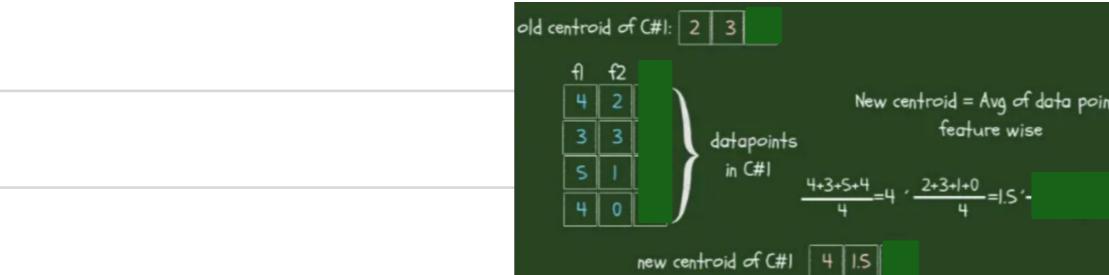


↳ 3. Update:

↳ Recalculate centroids as the means

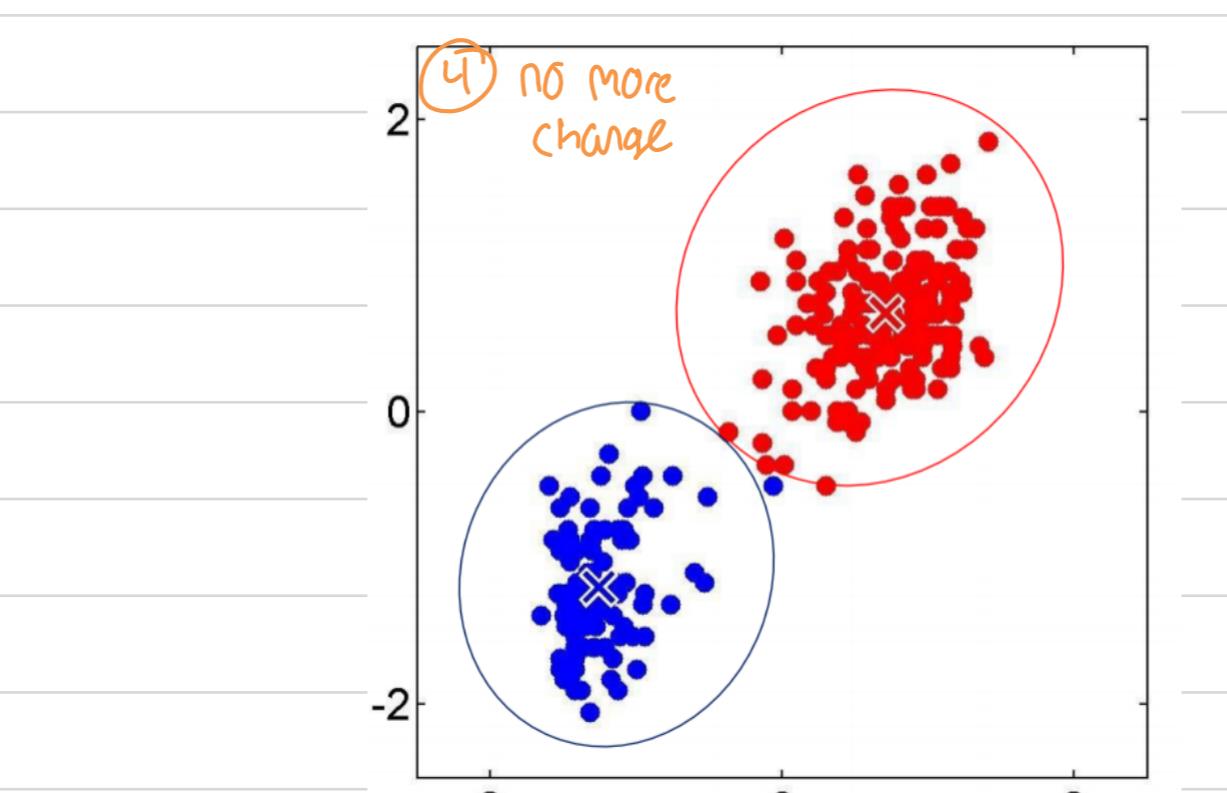
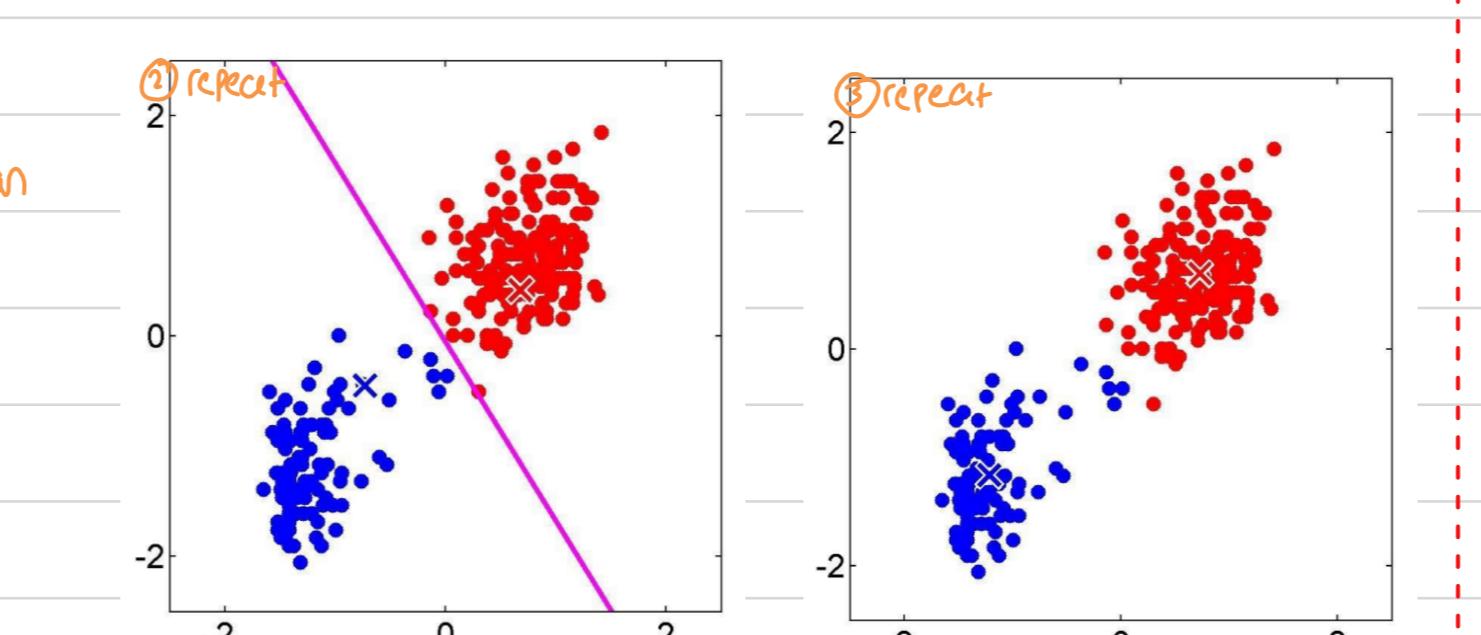
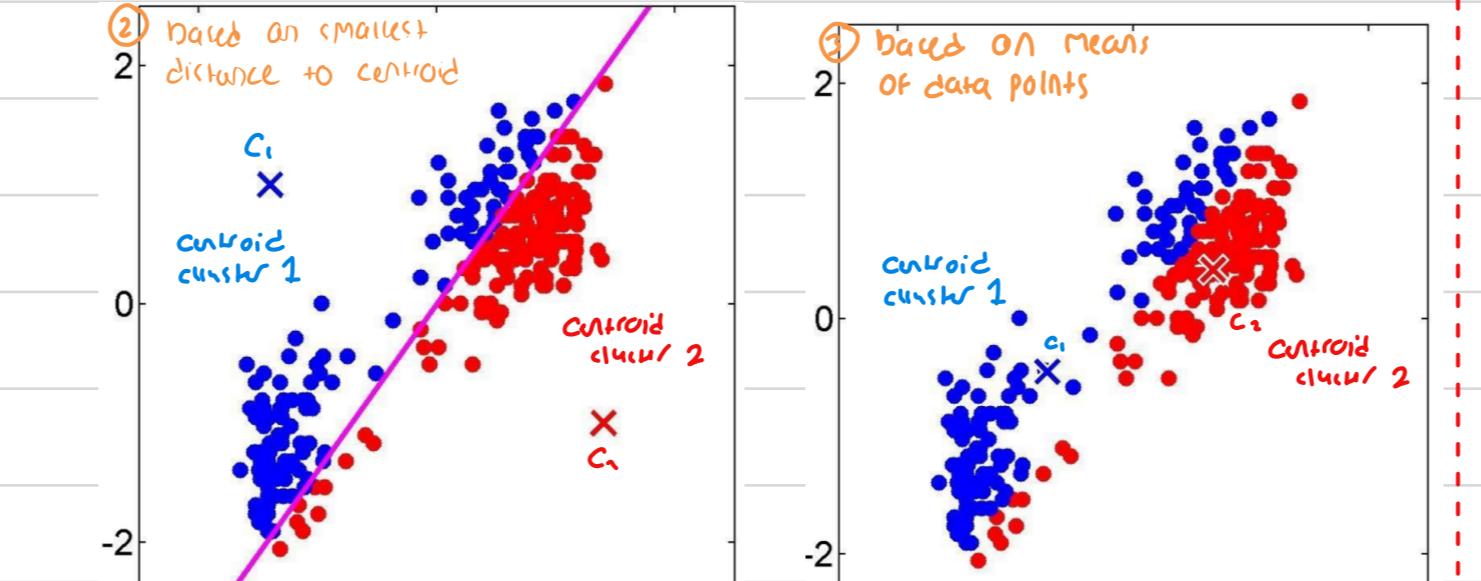
↳ distance of data point in the cluster, based on smallest distance to centroid

↳ by simple Mean ( $k$ -mean)



↳ 4. Convergence:

↳ Repeat Step 2 & 3 until the change in centroids below threshold / max iteration



NOT IN SCOPE

## Hierarchical Clustering

↳ recursively merge & divide clusters based on similarity

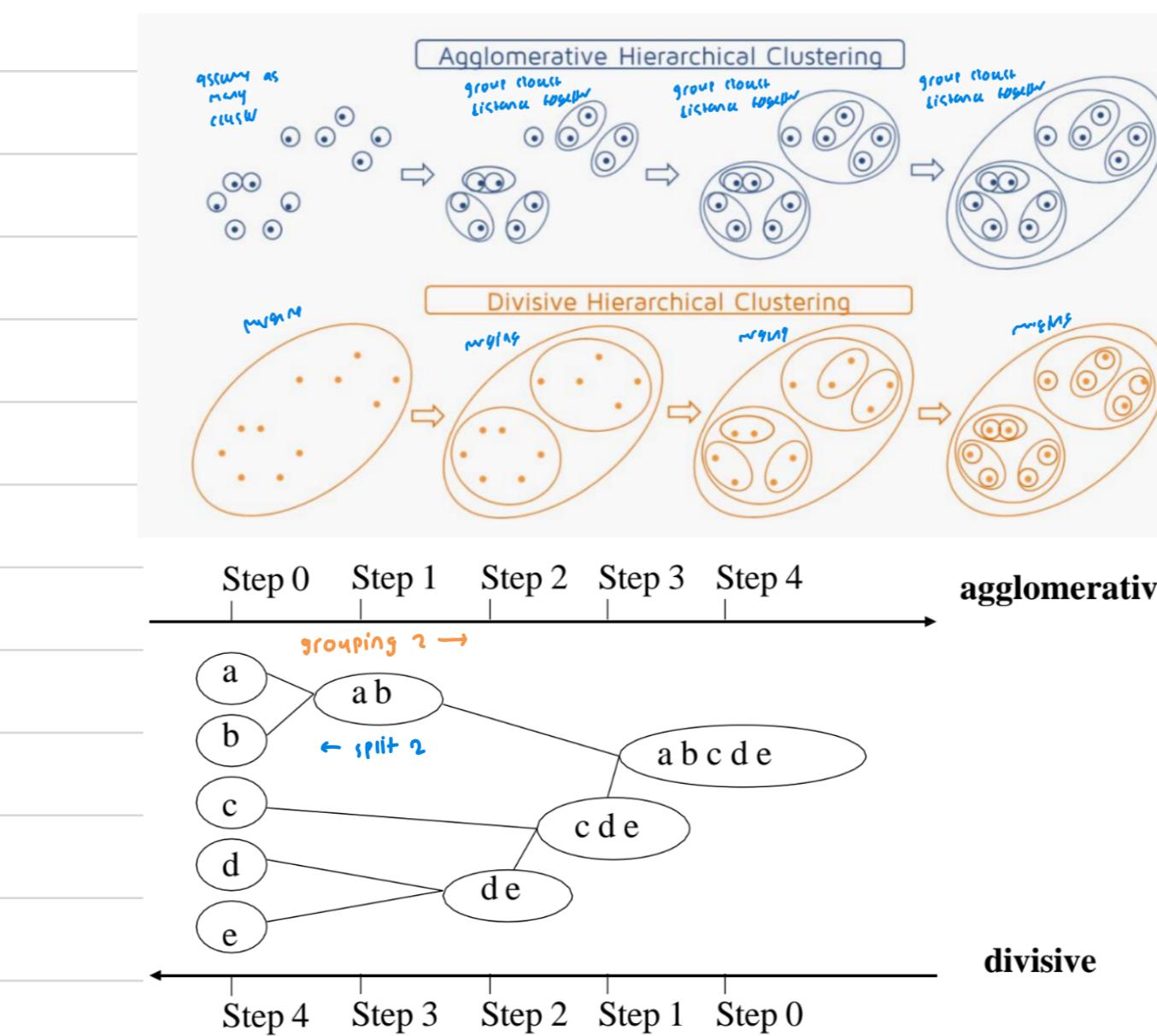
↳ 2 types:-

### Agglomerative (Bottom-up)

↳ recursively merge two clusters together w/ the lowest between-cluster dissimilarity

### Divisive (Top-down)

↳ Split the least coherent cluster



## How to calculate Dissimilarity :-

• Have to measure the dissimilarity for two disjoint groups  $G$  and  $H$ ,  $D(G, H)$  is computed from pairwise dissimilarities  $D(i, j)$  with  $i \in G, j \in H$

↳ Single Linkage: tends to yield extended clusters.

$$D_{SL}(G, H) = \min_{i \in G, j \in H} D(i, j)$$

↳ Complete Linkage: tends to yield round clusters.

$$D_{CL}(G, H) = \max_{i \in G, j \in H} D(i, j)$$

↳ Group Average: tradeoff between them. Not invariant under monotone increasing transform.

ave mean point

$$D_{GA}(G, H) = \frac{1}{n_G n_H} \sum_{i \in G, j \in H} D(i, j)$$

## Spectral Clustering

↳ datapoints as vertices of a Graph, all vertices are connected by an Edge, all edges have Weight

↳ weight ↑, adjacent V are very similar

↳ weight ↓, adjacent V are very dissimilar

↳ clustering using loss function

$$\text{cut}(A, B) = \sum_{u \in A, v \in B} w_{uv}$$

↳ Types of cut:-

• MinCut criterion

↳ Find partition  $(\bar{A}, \bar{B})$  that minimize  $\text{cut}(A, B)$

↳ ignored subgraph size result

• Normalized cut criterion

↳ More balanced partitions

↳ Steps:-

1. Compute a similarity graph

↳ i) Create an undirected graph

ii) Create adjacency matrix

$$W = \begin{bmatrix} 0 & w_{12} & w_{13} & \dots \\ w_{21} & 0 & w_{23} & \dots \\ w_{31} & w_{32} & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

2. Project data to low dimensional space

↳ transform the space so that when 2 point are close they are always in the same cluster & vice versa

↳ Compute Graph Laplacian to find clustering

Properties in Graph

$L = D - A$ , where  $A$  is the adjacency matrix and  $D$  is the degree matrix such that  $d_{ii} = \sum_j A_{ij}$ . Thus we get eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  where  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ .

1. If  $L$  has eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  where  $0 = \lambda_1 > \lambda_2 = \dots = \lambda_k$ , graph  $G$  has  $k$  connected components.

2. If  $G$  is connected, i.e.,  $\lambda_1 > \lambda_2 > \dots > \lambda_n$  is the algebraic connectivity of  $G$ . Thus, greater  $\lambda_2$  denotes the connectivity.

Computing Graph Laplacian

$L = D - A$ , where  $A$  is the adjacency matrix and  $D$  is the degree matrix such that  $d_{ii} = \sum_j A_{ij}$ . Thus we get eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  where  $0 = \lambda_1 > \lambda_2 = \dots = \lambda_k$ , graph  $G$  has  $k$  connected components.

Graph Laplacian for example above

$L = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 & 0 \\ 1 & 0 & 0 & 4 & 0 & 0 \\ 1 & 0 & 0 & 0 & 5 & 0 \\ 1 & 0 & 0 & 0 & 0 & 6 \end{bmatrix}$

$D = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 6 \end{bmatrix}$

$D^{-1/2} L D^{-1/2} = \begin{bmatrix} 1 & 0.5 & 0.33 & 0.25 & 0.2 & 0.17 \\ 0.5 & 2 & 0 & 0 & 0 & 0 \\ 0.33 & 0 & 3 & 0 & 0 & 0 \\ 0.25 & 0 & 0 & 4 & 0 & 0 \\ 0.2 & 0 & 0 & 0 & 5 & 0 \\ 0.17 & 0 & 0 & 0 & 0 & 6 \end{bmatrix}$

Graph Laplacian for example above

Graph Laplacian for example above